# Hybride & IA

LES NOUVELLES GÉNÉRATIONS DE MESURE FACE AUX ENJEUX DE LA MESURE MÉDIAS

Aurélie VANHEUVERZWYN Julien ROSANVALLON



# Avant-propos

L'écosystème média connaît actuellement une profonde transformation portée par la délinéarisation et la fragmentation croissantes des usages. Ce mouvement de fond se traduit par des nouveaux besoins pour l'ensemble des acteurs du marché et la mesure d'audience doit se transformer pour y répondre.

Les sources de données utiles aux mesures se multiplient: les données voie de retour - les data - couvrent désormais une large partie des usages et deviennent un complément intéressant aux données d'audience individuelles issues des panels pour la quantification robuste des audiences d'un plus grand nombre de contenus. Elles apportent en effet la précision et la granularité nécessaires à ce cas d'usage. Mais ces data véhiculent aussi des mythes qui laissent parfois penser que l'on pourrait se passer des panels et enquêtes. La réalité que nous tenterons d'éclairer, c'est que c'est bien la combinaison de ces sources de données qui permet la transformation des systèmes de mesure d'audience.

Le sujet n'est pas nouveau, et l'évolution vers des systèmes hybrides est déjà en marche dans de nombreux pays. Ce livre blanc a pour objectif d'analyser ces évolutions, d'un point de vue, historique, théorique, mais aussi au travers d'exemples concrets de mise en œuvre, et d'explorer les perspectives ouvertes par les développements récents de l'IA. Pour compléter cette lecture, plusieurs experts du secteur ont accepté de partager leur regard sur le sujet et ses enjeux.



# Remerciements

Nous tenons à remercier sincèrement toutes les personnes qui ont contribué à faire de ce livre blanc une réalité.

Un grand merci aux personnes qui ont contribué en apportant leurs témoignages sur ce sujet: Jean-Luc Chetrit (UDM/WFA), Koenraad Deridder (CIM), Pete Doe (Nielsen), Khaled El Serafy (BARB), Derrick Gray (Numeris), Mirko Marr (Mediapulse AG), Andrea Mezzasalma (dataBreeders), Valérie Morrisson (CESP), Zuber Nosimohomed (Kantar Media), Mario Paic (Ipsos).

Nous tenons aussi à remercier les équipes de Médiamétrie et notamment Laurence David, Laurence Deléchapt, Mathilde Jehan, Emmanuelle Le Goff, Marie Liutkus et Arnaud Philippe pour leurs précieuses contributions.



63

# Sommaire

**5** Synthèse et prochaines étapes

Préface	5
1 Introduction	9
<ul> <li>Une brève histoire de la mesure d'audience des médias</li> </ul>	10
L'âge de la donnée : une révolution numérique	12
Mythes et réalité d'un monde de mesure où la data remplace enquêtes et panels	14
Premières formes de régulation de la data	17
La nécessaire synthèse des approches	18
2 La réconciliation des data et des échantillons : les approches hybrides	20
Les bénéfices de l'hybridation	22
Principes & typologie de mesure hybrides	24
Les défis inhérents aux mesures hybrides	27
■ Echantillonner la data	29
Les exemples de mesure hybride dans le monde	30
<b>3</b> L'intelligence artificielle : de nouvelles perspectives pour les mesures	37
■ Une nouvelle ère pour la mesure d'audience ?	38
■ IA et données synthétiques	40
■ Enjeux éthiques, transparence et explicabilité	43
■ IA et accès aux données d'audience	44
4 L'hybride & l'IA vus par les experts	45

Hybride & IA Préface

### **Préface**

# Entre fragmentation des usages et puissance algorithmique : la mesure en mutation

La transformation des systèmes de mesure d'audience est aujourd'hui portée par deux dynamiques majeures : l'hybridation des sources de données et l'intégration de l'intelligence artificielle. Ces évolutions ne sont pas simplement techniques ; elles redéfinissent les fondements mêmes de la mesure, sa légitimité, sa gouvernance et sa capacité à accompagner les mutations des usages médias.

L'hybridation, d'abord, répond à une exigence de représentativité et de granularité. Les panels, socle historique de la mesure, offrent une rigueur méthodologique et une stabilité précieuse. Mais ils peinent à capter l'infinie diversité des comportements numériques, fragmentés, multiécrans, multi-contextes. L'apport des données massives – logs, impressions, navigation – permet d'élargir le spectre, de détecter les signaux faibles, de mieux comprendre les parcours. Encore fautil que cette hybridation repose sur des principes clairs:complémentarité des sources, transparence des modèles, auditabilité des résultats.

C'est là que l'IA entre en scène. Elle permet de modéliser, de redresser, de prédire. Elle peut enrichir les données, détecter les anomalies, optimiser les traitements. Mais elle soulève aussi des questions éthiques majeures : comment garantir l'absence de biais ? Comment expliquer les décisions algorithmiques ? Comment assurer la souveraineté des acteurs sur leurs données ? L'IA ne doit pas être une boîte noire, mais un levier de confiance et de progrès.

« L'IA ne doit pas être une boîte noire, mais un levier de confiance et de progrès » Dans ce contexte, la mesure est plus que jamais un enjeu stratégique pour les marques et les médias. Elle ne peut plus être pensée comme un outil figé, mais comme un système vivant, évolutif, co-construit. La mesure cross médias, à laquelle nous travaillons tous ensemble, incarne cette exigence. Elle vise à restituer une vision dédupliquée de l'exposition aux campagnes, tous canaux confondus – aujourd'hui TV, digital, CTV, plateformes vidéo et demain aux autres médias. Pour les marques, il s'agit de mieux piloter la pression publicitaire, d'éviter les surexpositions et de renforcer l'efficacité des dispositifs.

À l'international, les pilotes de l'initiative de la WFA (World Federation of Advertisers) comme Origin (UK) ou Aquila (USA), sont portés par les annonceurs. Ils inspirent les développements français menés par Médiamétrie avec le soutien de l'Union des marques. D'ores et déjà des tests sur les campagnes TV+CTV, combinant watermarking, logs adservers et modélisation ont permis d'identifier des gains de performance significatifs. Ces résultats confirment que l'hybridation des approches, soutenue par l'IA, permet de mieux refléter la réalité des usages et d'optimiser les stratégies médias.

L'intelligence artificielle, en enrichissant les modèles et en optimisant les traitements, permet d'accélérer cette transformation, à condition de rester transparente et responsable. Elle doit être au service de la comparabilité, de la performance collective et de la souveraineté des acteurs.

La mesure cross médias, augmentée par l'IA, est plus que jamais une opportunité pour les marques : celle de construire des dispositifs plus justes, plus inclusifs, plus efficaces. Mais c'est aussi une responsabilité collective : celle de garantir que la technologie reste au service de la transparence, de la comparabilité et de la confiance.

Jean-Luc Chetrit Directeur Général de l'UDM Hybride & IA Préface

### **Préface**

# Pourquoi l'avenir de la mesure d'audience doit rester hybride

C'est peu de dire que la consommation médias des Français s'est fragmentée au cours des dernières décennies. Avant le lancement de la TNT, les Français avaient le choix entre six chaînes de télévision, aujourd'hui plus de 1000 chaînes sont disponibles. L'offreaudios'estégalement beaucoup enrichie. À cette profusion de l'offre, s'ajoute la multiplication des modes de consommation dans un monde de l'hyperdistribution où les contenus sont disponibles au travers d'offres et de canaux multiples, notamment via les plateformes OTT et les réseaux sociaux, particulièrement importants pour les plus jeunes. Dans ce contexte, la mesure d'audience court après les usages et un panel, même de taille significative, ne peut plus rendre compte de la multiplicité de ces usages.

Augmenter la taille des échantillons trouve ses limites dans la difficulté croissante à joindre les personnes, notamment les plus jeunes, et dans les coûts additionnels engendrés. La bonne nouvelle pour les mesureurs est que tous nos actes numériques laissent des traces, offrant une vaste quantité de données à disposition (connexions, navigations, streams...) et de nouvelles opportunités pour les mesures hybrides.

Pour le CESP, tiers de confiance en charge de l'audit des études d'audience, ces mesures hybrides requièrent de s'assurer de la qualité et donc de la fiabilité des résultats obtenus, sur plusieurs plans. Tout d'abord, la qualité de la donnée source est essentielle (« garbage in, garbage out » disent les Anglais!). Cela couvre à la fois les panels ou les enquêtes mais aussi les données dites passives (server-side logs, device-logs...), avec deux dimensions: la certification des volumes de données mais aussi la qualité de ces données. Ensuite, le sujet du rapprochement de ces sources se pose avec de multiples techniques d'hybridation possibles, clairement exposées

dans ce livre blanc. Toutes les mesures d'audience comportent aujourd'hui au moins un modèle, voire de multiples modèles qui s'enchaînent, ce qui pose deux grandes questions. Premièrement quelle est la source de vérité qui permet de valider la qualité du modèle ? Deuxièmement, quelles sont la robustesse et la précision de la donnée en sortie ? Nous observons que les calages ou calibrations se font bien souvent à un niveau macro, mais dans le but de produire des résultats micro, c'est-à-dire au niveau de cibles, en réponse aux attentes des utilisateurs, ce qui nécessite une grande vigilance dans l'utilisation des résultats au niveau opérationnel.

Demain peut-on imaginer un monde de la mesure avec uniquement de la donnée passive matinée d'IA générative qui viendra la compléter et l'augmenter ? Je ne le crois pas.

évolutions actuelles renforcent Les trois ďun contrôlé avantages panel géré centralement. Premièrement échapper phénomène des bots qui, enrichis de l'IA, pourront générer du trafic invalide de plus en plus sophistiqué dans les signaux digitaux. Deuxièmement, avec un cadre réglementaire de plus en plus restreint, exigeant un consentement toujours plus explicite, un panel fondé sur le consentement éclairé des participants est plus résilient. Enfin un panel piloté et financé de manière multipartite évite les phénomènes d'automesure et permet au marché de partager une même currency pour le bénéfice des annonceurs.

> « Le futur de la mesure d'audience sera nécessairement hybride »

Les 10 dernières années passées au CESP, pendant lesquelles nous avons audité chaque année toutes les mesures d'audience du marché français, mais aussi d'autres pays, en Europe, en Asie, en Afrique et au Moyen-Orient, m'ont convaincue qu'aucune méthodologie n'est parfaite. Enquête ou données passives, chacune a ses avantages et ses inconvénients et les panels single source, potentiel graal, trouvent eux-mêmes leurs limites dans leur coût et leur acceptabilité par les panélistes qui n'ont pas nécessairement envie d'être suivis dans l'ensemble de leurs usages médias.

Ma conviction est donc que le futur de la mesure d'audience sera nécessairement hybride et que les panels et les enquêtes garderont un rôle central pour calibrer les modèles et étiqueter les données de manière fiable afin d'alimenter les hybridations de demain.

Valérie Morrisson Directrice Générale du CESP

# Une brève histoire de la mesure d'audience des médias

La mesure d'audience des médias est une pratique aussi ancienne que les médias eux-mêmes. Elle naît du besoin fondamental de connaître son public : qui regarde, qui écoute, qui lit ?

Pour le média Presse, les données de diffusion des titres constituaient une information essentielle qui renseignait les journaux sur l'appétence du public. Cette information ne constitue pas une audience, puisqu'un exemplaire diffusé peut être vu par plusieurs personnes. Mais c'est une information sur les goûts du public (et sa capacité à acheter un numéro, à s'abonner...). Avant qu'une véritable mesure d'audience de la presse soit lancée, cette donnée était souvent complétée par le courrier des lecteurs apportant un prisme qualitatif. Touchez à une rubrique que le public affectionne et il le fera savoir à la rédaction dans un courrier abondant.

Pour les médias audiovisuels, la situation était plus complexe puisqu'au-delà de données du courrier des auditeurs / téléspectateurs, la diffusion par ondes « hertziennes » rendait impossible la production de données de diffusion similaires à celles de la presse : on connait le nombre d'émetteurs physiques, mais aucune information disponible sur les récepteurs (le public). C'est ce qui a encouragé le développement rapide des premières mesures d'audience des médias radio et télévision, rapidement après leur développement auprès du grand public.

Cette connaissance du public, essentielle pour les éditeurs comme pour les annonceurs, a ainsi toujours été au cœur de la valorisation des contenus et des espaces publicitaires. Pourtant, les outils et les méthodes ont profondément évolué au fil des décennies, à mesure que les usages et les technologies se sont transformés.

### Des premiers pas : l'ère du déclaratif et des enquêtes téléphoniques

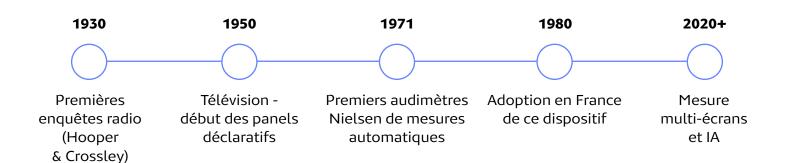
Les premières mesures d'audience remontent aux années 1930, à l'époque de la radio. Aux États-Unis, des instituts pionniers comme Hooper et Crossley recouraient à des enquêtes téléphoniques pour estimer l'écoute des programmes. Quelques années plus tard, en 1949, Nielsen acquit les droits de l'Audimètre, une nouvelle technologie reliée à un émetteur radio qui permettait de capturer ce que les gens écoutaient.

Avec l'arrivée de la télévision dans les foyers dans les années 1950, le besoin d'une mesure plus précise et plus continue s'est fait sentir. En France, les premières mesures reposaient alors sur une méthodologie proche de celle de la radio. Mais à la fin des années 60, des systèmes de mesure automatique se développèrent (le boîtier installé auprès d'un panel de téléspectateurs). En 1971, Nielsen a lancé un audimètre capable de stocker les données d'utilisation des téléviseurs et de les transmettre pendant la nuit via une ligne téléphonique. Le dispositif fut adopté en France au début des années 80. Les technologies de mesure ont ensuite évolué, s'adaptant aux transformations des modes de diffusion et de distribution des médias audiovisuels, afin de permettre la restitution des audiences dans toutes les configurations possibles (tous lieux, tous écrans, toutes temporalités). Les nouvelles technologies de mesure comprennent désormais des logiciels installés sur des appareils numériques (Realitymine, Ipsos Mediacell, Gemius, etc.) ou des appareils individuels portables sous forme de bracelets connectés ou de montres (Médiamétrie RateOnAir, Nielsen PPM, Gfk Mediawatch).

# Echantillons et panels au cœur de la mesure

Toutes ces mesures reposent sur les données d'échantillons et de panels qui deviennent des standards des mesures d'audience. Cela s'explique tant par la contrainte pratique de ne pas pouvoir interroger tout le monde, que de la nécessité de faire reposer la mesure sur une donnée individuelle qui représente les usages et consommations médias de personnes.

La promesse est forte : un petit nombre de foyers, soigneusement choisis, doit représenter les comportements de millions de personnes. Cette approche fondée sur la théorie des sondages allie rigueur scientifique, transparence méthodologique et gouvernance partagée. Elle permet la production d'études qui deviennent de véritables monnaies et un langage commun pour le marché des médias.



# L'âge de la donnée : une révolution numérique

### « Graal de la mesure : allier représentativité et granularité infinie »

Avec l'essor d'internet grand public, du digital et des appareils connectés dans les années 2000, un nouvel horizon s'ouvre : celui de nouvelles données massives qui enregistrent toutes les actions de consommations numériques. On touche alors au Graal de la mesure : allier représentativité et granularité infinie. Le développement du web s'est ainsi accompagné de mesures analytics qui comptabilisent chaque clic des internautes. Une mesure qui n'a eu de cesse de se développer et contribua au succès publicitaire du média digital en permettant une collecte massive de données.

Toutes ces nouvelles données voie de retour et outils de mesure qui les exploitent deviennent la référence dans l'écosystème numérique. En particulier, le marché publicitaire digital adopta largement le concept d'impressions publicitaires. Une donnée qui mesure la pression publicitaire. Un concept de pression qui ressemble à celui de GRP (le nombre de contacts sur cibles). Un point majeur distingue les deux. Le GRP est en effet produit sur cible et nécessite donc de savoir qui a vu la publicité.

Le développement de cette économie de la data fut spectaculaire et permis à quelques acteurs, aux écosystèmes et services uniques (de recherche, de réseaux sociaux...), de développer un modèle publicitaire puissant.

# Les promesses et bénéfices de la data face aux échantillons

L'arrivée du numérique a bouleversé l'économie des médias. Avec elle, une nouvelle ressource est devenue centrale : la data, donnée comportementale collectée de façon massive. À la différence des méthodes traditionnelles fondées sur des panels ou des enquêtes, la data permettait, en théorie, de capter chaque action : une page vue, un clic, une vidéo démarrée, un contenu partagé. Une révolution s'annonce : la fin de la mesure par sondage au profit d'une mesure exhaustive, en temps réel et granulaire.

La production de résultats d'audience sur la base d'échantillons est soumise à un ensemble de contraintes bien connues. Les data promettent de se libérer de ces contraintes :

### - Les échantillons sont soumis à des biais de sélection :

À l'exception de la Statistique Publique qui dispose d'une base de sondage de l'ensemble de la population, les échantillons ne sont jamais constitués de manière totalement aléatoire. La méthode des quotas, qui consiste à aligner la structure des panels et des échantillons sur celle de la population, repose sur l'hypothèse selon laquelle le fait de répondre à l'enquête dépend de caractéristiques socio-démographiques et non des comportements d'audience. Mais cette hypothèse ne peut jamais être totalement vérifiée.

### - Les échantillons sont soumis à des limites de précision :

Le principe d'un sondage est d'estimer des informations sur une population en n'interrogeant qu'une partie de celle-ci, l'enjeu étant de sélectionner un échantillon représentatif de la

« Face aux limites des échantillons, connues et documentées, les data ouvrent un nouveau champ d'exploration et de nouvelles potentialités pour les mesures d'audience des médias. »

population. Le rapport entre la taille de la population et la taille de l'échantillon permet de mesurer le nombre de personnes que chaque interviewé / panéliste représente. Pour la mesure de la télévision, un panéliste représente en moyenne environ 6000 Français. On comprend ainsi aisément les limites des échantillons pour la mesure d'usages très fragmentés.

### - Les résultats produits à partir d'échantillons sont soumis à l'erreur d'échantillonnage :

L'erreur d'échantillonnage est due à la variabilité inhérente au processus de sélection d'un échantillon. Deux échantillons de même taille sélectionnés en suivant strictement le même protocole donneront des résultats différents.

Cette variabilité est documentée scientifiquement et se quantifie par des indicateurs tels que l'intervalle de confiance.

Les sociétés qui opèrent les enquêtes et les recrutements de panels mettent en œuvre des méthodes rigoureuses afin de limiter la portée de ces contraintes, sous le contrôle des tiers qui les auditent. En France, ces audits externes sont menés par le CESP. Un rapport très détaillé pour chaque média analyse les caractéristiques méthodologiques de l'étude et la qualité de la réalisation.

Face aux limites des échantillons, connues et documentées, les data ouvrent un nouveau champ d'exploration et de nouvelles potentialités pour les mesures d'audience des médias.



# Mythes et réalité d'un monde de mesure où la data remplace enquêtes et panels

Face à ces perspectives alléchantes, des discours provocateurs ont alors fleuri : « le panel est mort », « la data va tout révéler », « plus besoin d'échantillons, on a la totalité ». Dans l'industrie, certains annonceurs et agences ont vu dans la data une opportunité de reprendre le contrôle sur la chaîne de valeur média. Du côté des plateformes, la maîtrise de ces données devenait un levier de puissance : qui possède la data, possède la mesure, donc la valeur. Enfin, pour certains éditeurs, exhaustivité signifiait des audiences supplémentaires. au-delà volatilité. d'une baisse la

Mais cette promesse s'est aussi heurtée à une réalité plus complexe. Trois mythes illustrent la complexité de cette nouvelle réalité de ce monde de data.

# Mythe 1: « La data est exhaustive et offre une vision globale »

Les données sont collectées à de nombreux niveaux: des boxes, des applications, des téléviseurs connectés, des terminaux mobiles, des plateformes.... Chaque jeu de données collecté sur un de ces segments apporte une vue exhaustive de cet univers. La plateforme A a accès à un volume très important de données sur ses utilisateurs, mais ne sait rien de l'usage des utilisateurs lorsqu'ils évoluent sur la plateforme B. Les données sont très riches, mais elles ne sont le reflet que du segment de population qu'elles analysent. Cette limitation est importante si on veut utiliser ces data pour avoir une représentation des usages de l'ensemble de la population.

Par ailleurs, au-delà de ce sujet de couverture de l'ensemble des usages, la data étant propriétaire, elle est souvent produite selon des règles de calculs et des normes spécifiques. Ceux qui la possèdent imposent souvent leurs propres métriques, leurs propres standards, leurs propres modèles. L'absence de normalisation rend complexe les comparaisons de différentes data.

Comprenant cette limite, les institutions de l'Union Européenne ont pris récemment une position claire dans le cadre du Règlement Européen sur la Liberté des Médias. Ainsi, l'article 24 de ce règlement promeut-il l'application des principes de mesure d'audience marché, incluant explicitement celui de la comparabilité. A ces principes s'ajoute un cadre strict, en matière de protection des données personnelles des utilisateurs de service numérique qui a un impact significatif sur ce que les outils analytics peuvent produire. Le RGPD notamment, mais aussi la Directive ePrivacy, posent le principe d'un traitement des données analytics sur la base du consentement des utilisateurs pour la mesure de l'audience. Des exceptions apparaissent toutefois au niveau de certains Etats de l'Union Européenne. C'est notamment le cas de la France où la Cnil a publié des Lignes Directrices et une Recommandation permettant la collecte de données sous l'égide d'une exemption de consentement. Le traitement des données collectées est alors strictement encadré permettant un comptage simple sur des éléments de volume intéressants mais limitant la profondeur des informations.

Aux États-Unis, lors de ses travaux pour passer à son nouveau système TV+Big Data, Nielsen a été amené à comparer les données de Smart TV avec celles de leur panel. La comparaison a pu être faite au niveau de chaque appareil mesuré par les deux systèmes. Il en est ressorti que les data des Smart TV avaient un périmètre plus restreint en termes de chaînes et de programmes que celles des mesures de Nielsen.

Cela s'explique par le fait que les données issues des Smart TV sont généralement recueillies par le biais d'une technologie de fingerprinting video ou audio. Un logiciel embarqué dans le système d'exploitation téléviseur identifie (une fois le consentement de l'utilisateur recueilli) les contenus et chaînes diffusés. Le problème est que ce système repose sur la capacité de ces acteurs à créer une bibliothèque de ces contenus et chaînes. Ces bibliothèques sont importantes, mais ne peuvent être exhaustives, soit pour des raisons pratiques, mais aussi parfois pour des raisons contractuelles.

De façon alternative, les data des téléviseurs peuvent être collectées par le biais des boxes. Dans ce cas, la collecte de données (toujours soumise au recueil préalable du consentement) sera plus exhaustive. Mais son périmètre est lui aussi limité aux services de TV/Vidéo contrôlés directement par cet opérateur. Dans les faits ces données ne capturent également qu'une partie des usages TV/Vidéo des utilisateurs.

### Mythe 2: « La data est exacte et infaillible »

Sur un plan théorique, au-delà des limites liées à ce que les data ne savent pas mesurer, celles-ci sont exhaustives. Elles ne sont donc pas entachées d'erreur d'échantillonnage. Néanmoins, la mise en œuvre opérationnelle de ces mesures peut induire des erreurs de différentes natures :

### - Les limites techniques inhérentes à la mesure :

Certains utilisateurs des outils analytics exploitent des indicateurs sans en connaître les limites. Les outils analytics ne mesurent pas des individus. Ils mesurent au mieux un nombre de terminaux, voire plus souvent un nombre de navigateurs, ou de cookies. Ces données sont donc une estimation globalement fausse du nombre de personnes qui visitent un site pour un ensemble de raisons :

Fragmentation des usages sur plusieurs terminaux: Les utilisateurs ont plusieurs terminaux. Or chaque terminal est souvent identifié différemment par les mesures qui ne savent pas bien les lier.

Effacement des cookies et surf anonyme : La plupart des terminaux permettent d'effacer simplement les cookies ou de surfer en utilisant un mode qui limite les données personnelles et de navigation transmises aux sites.

Partage de terminaux: Certains écrans sont partagés par plusieurs utilisateurs, rendant difficile l'identification des individus.

### - Des erreurs liées au paramétrage :

Internet est un réseau interconnecté sur lequel les hommes ne sont pas les seuls à « naviguer ». Les bots et programmes informatiques des acteurs du numérique parcourent les contenus des pages pour les référencer, les analyser, les comprendre... Un trafic colossal qui peut fausser les mesures analytics si celles-ci ne sont pas correctement paramétrées. L'IAB/ABC édite notamment une liste des bots et spiders à exclure des mesures. Cette liste est mise à jour en permanence. Les principaux outils excluent donc bien ce trafic.

Mais il y a toujours une possibilité qu'un de ces bots créés ne soit pas déclaré et génère un trafic artificiel (même si d'autres méthodes existent pour les identifier).

#### - La fraude:

La fraude est un sujet important qui peut permettre à des tiers malveillants de générer des recettes publicitaires en simulant un trafic artificiel. Ainsi les pirates à l'origine du Hyphbot en 2017 ont créé un système qui générait de faux sites Premium identifiés comme tels par les réseaux publicitaires. Ces faux sites étaient ensuite alimentés par un trafic de robots qui généraient des millions d'impressions fictives. La fraude fut détectée par un outil de détection de la fraude (AdForm). De nombreux acteurs proposent d'ailleurs aujourd'hui leurs services pour détecter et stopper ces tentatives.

# Mythe 3 : « La data remplacera les panels »

Ce mythe, qui fait la synthèse de deux autres, repose sur l'idée que les données massives suffisent à décrire fidèlement les comportements. Or un ensemble de données, même exhaustif, n'a de valeur que s'il couvre un univers pertinent

et qu'il permet de calculer des indicateurs intelligibles pour le marché. La mesure d'audience, apporte une réponse à plusieurs dimensions clés qui ne sont pas nativement couvertes par les données voie de retour:

### - Individuelle:

En suivant l'audience d'individus et non de machines.

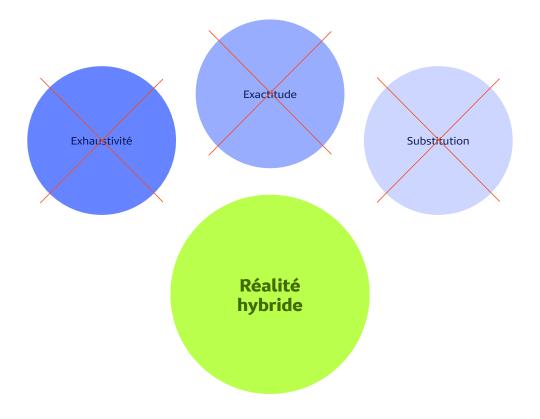
### - Transversale:

En permettant de dédupliquer les usages des individus dont les data sont potentiellement dans des silos étanches (les data d'une plateforme et d'une autre par exemple).

### - Longitudinale:

En suivant les comportements des mêmes individus dans le temps.

La data, elle, a une richesse infinie, mais manque de ces qualités essentielles à la production d'une mesure d'audience. L'un ne peut donc remplacer l'autre. Ce sont deux outils différents qui répondent à des logiques distinctes, et qui peuvent - ensemble - enrichir la compréhension des usages.



### Premières formes de **régulation de la data**

Tous ces facteurs qui viennent limiter ou perturber la mesure sont autant de raisons qui ont vu émerger au fil des années :

- Des normes, notamment dans le champ publicitaire, comme celle de l'IAB (Internet Advertising Bureau), le MRC (Media Ratings Council). L'IAB a publié plusieurs rapports et guidelines, comme le « Interactive Audience Measurement and Advertising Campaign Reporting and Audit Guidelines », ou encore les « Digital Video Ad Impression Measurement Guidelines ».
- Des audits & certifications des données analytics éditoriales (comme ceux réalisés par l'ACPM en France, ou les ABC dans le monde). Aux États-Unis le Media Ratings Council propose ses services pour auditer des solutions de mesure (panel et analytics). En France, le CESP (Centre d'Études des Supports Publicitaires), en complément de l'étude régulière des mesures d'audience de référence des médias, réalise des audits spécifiques comme celui des outils de mesure de la visibilité.

- Des acteurs spécialisés dans l'analyse de la fraude (Human, Integral Ad Science, AdLoox, DoubleVerify...) qui proposent leurs services pour apporter une analyse poussée, incluant la détection de la fraude, dans les métriques publicitaires.
- Si l'on met de côté les contrôles opérés par les acteurs comme l'ACPM ou les ABC, la plupart de ces outils de contrôles restent optionnels et leur utilisation repose sur le choix des acteurs.

# La nécessaire synthèse des approches

L'univers des médias est de plus en plus complexe. Ses usages se segmentent et se fragmentent. En France, le lancement de la TNT en 2005 a été autant un moteur pour le développement des usages du média, en développant l'offre, qu'il a reconfiguré un paysage audiovisuel qui reposait sur quelques chaînes. Dans l'univers du digital, ces effets de fragmentation se sont encore accélérés. Si moins d'une dizaine de marques digitales réunissent plus de 10 millions d'utilisateurs chaque jour, près de 900 sont visitées par plus d'un million de personnes chaque mois. Dans l'univers de l'audio, l'émergence des podcasts fait exploser l'offre disponible à un instant donné, et donc contribue à la fragmentation des usages. À titre d'illustration, sur une plateforme comme Spotify, plus de 6 millions de podcasts sont disponibles (source Spotify - Décembre 2024). L'Arcom de son côté identifie qu'en France, plus de 10 millions d'épisodes de podcasts sont disponibles. Ces chiffres attestent des profondes transformations du paysage et de ses usages qui sont potentiellement ouverts par cette nouvelle offre.

Face à cette réalité, les limites des échantillons décrites plus haut se font ressentir pour restituer une certaine granularité des usages. C'est pourquoi certaines mesures ont adapté la taille des échantillons utilisés. Ce fut le cas notamment de la mesure de la télévision qui fit largement progresser son échantillon au fil des années (2300 foyers en 1999 – 5500 foyers et un panel complémentaire de 5500 individus en 2025) pour accompagner ces transformations. Nous observons des changements similaires dans la plupart des pays.

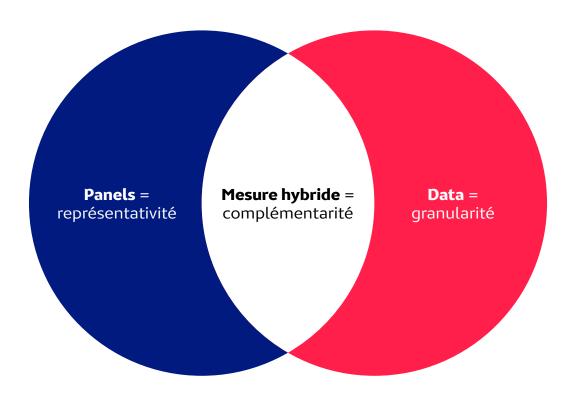
Dès la fin des années 90 et le début des années 2000, le digital et la fragmentation qui l'accompagne voient se développer des outils analytics qui tirent parti des données de voie de retour. Ces nouvelles données évoluent alors dans des silos, les utilisateurs n'ont accès qu'à leurs seules données. Des données qui offrent néanmoins une précision jusque-là inconnue des usages. En parallèle, des mesures reposant sur des panels se développent (JupiterMMXI, NetRatings, Netvalue...). Les comparaisons entre les données issues de ces deux outils se multiplient, créant de nombreux débats car les mesures analytics mesurent généralement deux à trois fois plus d'utilisateurs que celles reposant sur des panels. La différence s'explique notamment par l'effacement des cookies qui sur-estime le nombre de visiteurs (puisqu'après chaque effacement de cookies, le système considère qu'il est en présence d'un nouvel utilisateur).

Le développement des mesures basées sur les data et leur comparaison avec celles issues des panels se pose aussi rapidement pour la télévision. Les opérateurs disposent en effet en temps réel du nombre de boxes connectées à une chaîne à un instant donné, ce qui permit le développement de mesures alternatives à la mesure par panel.

# Les premières mesures TV issues des données des boxes triple play

Ainsi lors du déploiement des boxes tripleplay des opérateurs au début des années 2000, des observateurs prédisaient la fin de la mesure d'audience de la télévision telle qu'on la connaissait. Les données de boxes seraient la nouvelle référence. Médiamétrie développa pendant quelques années une mesure de la télévision basée sur des données issues des voies de retour (Digitime). Mais la réalité opérationnelle plus complexe démontra que ces nouvelles données, si elles apportaient une richesse additionnelle, et notamment la possibilité d'analyser des scores en temps réel, ne pourraient remplacer celles du Médiamat. Une mesure d'audience traduit les comportements individuels du public. Et pas seulement la connexion de boxes à une chaîne. L'écran est-il allumé? Qui est présent devant l'écran? Comment faire le lien entre les usages d'un individu sur différentes boxes (sa TV à domicile, sa consommation lors d'un match dans un bar, ou encore d'émissions en replay sur son mobile)? Par ailleurs, si ces nouvelles données reflètent les usages d'un nombre croissant et majoritaire d'utilisateurs, dans de nombreux pays les modes de réception analogique sans voie de retour restent importants (En France, en 2025, un tiers de la population accède à la télévision par réception hertzienne sur un poste. Pour 13% d'entre eux, c'est le mode de réception exclusif). Or le profil de consommation de ces utilisateurs est très différent des utilisateurs « connectés ». Autant de questions centrales auxquelles une donnée de box seule ne peut répondre, et dont les réponses sont fournies par les mesures historiques par panel.

Cet exemple illustre clairement la nécessité de bâtir des mesures qui tireront partie des panels et des données voie de retour.



# La réconciliation des data et des échantillons: les approches hybrides

# La réconciliation des data et des échantillons : les approches hybrides

Derrière une opposition apparente entre les mesures centrées sur les individus et celles data, le constat progressivement : ni les panels seuls, ni les données seules ne suffisent à décrire fidèlement la réalité des usages de demain. Les données panels sont pertinentes dans un monde moins fragmenté. L'émergence des plateformes et des consommations à la demande font exploser la fragmentation des usages et nécessite de faire évoluer les approches. Symétriquement les limites des data, particulièrement dans un cadre européen très protecteur des données personnelles des utilisateurs, sont avérées. La solution est de combiner les données. Marier panels et data pour créer une nouvelle mesure plus granulaire et plus précise.

C'est le début de l'ère des mesures **hybrides**.

### Un principe ancien

Hybrider, c'est marier plusieurs sources de données d'origine différente pour en créer une nouvelle plus riche ou plus précise.

Le concept d'hybridation n'est pas récent. L'utilisation d'information auxiliaire est classique en théorie des sondages: le redressement d'un échantillon sur les données provenant du recensement de la population en est une parfaite illustration. Les exemples d'hybridation de différentes sources de données dans l'univers des médias sont nombreux.

La nouveauté est l'apparition d'une nouvelle génération de mesure hybride qui marie des données individuelles, « people data », avec des données voie de retour exhaustives « device / app data ». Une génération de mesure qui promet de pallier tant les défauts des données de panels et d'enquêtes que de celles des données exhaustives.

La mise en œuvre de ces nouvelles mesures suppose de nouveaux outils, de nouvelles compétences, et souvent des modèles statistiques complexes.

# Une histoire faite de recherche d'équilibre et de réinventions

« Chaque révolution technologique appelle une réinvention méthodologique »

L'histoire de la mesure d'audience est une histoire de recherche d'équilibre permanent : entre représentativité et exhaustivité, entre transparence et propriété des données, entre standards collectifs et besoins particuliers. Mais c'est aussi une histoire d'innovation permanente, où chaque révolution technologique appelle une réinvention méthodologique.

Aujourd'hui, face à la fragmentation des usages, au développement des plateformes et à la multiplication des données, la mesure d'audience est plus que jamais un enjeu stratégique. C'est dans cette continuité que s'inscrit la nécessité d'une mesure hybride, ouverte, gouvernée et intelligible, au service de tout l'écosystème média.

### Les bénéfices de l'hybridation

Les approches hybrides ont de nombreux atouts qui dépassent largement le seul accroissement de la précision des résultats.

Les médias numériques se caractérisent par une consommation toujours plus fragmentée : entre les écrans, les plateformes, les formats et les moments de consommation, il est de plus en plus difficile de mesurer tous les aspects d'un média avec un seul et même dispositif.

C'est dans ce contexte qu'émerge la **mesure hybride**, non comme une méthode unique, mais comme une **famille d'approches** qui visent à réconcilier des sources de données hétérogènes pour reconstituer une image plus fidèle des comportements médias.

La mesure hybride permet de tirer parti des forces complémentaires des différents dispositifs : panels, enquêtes, données numériques, CRM, logs, métadonnées techniques... Encore faut-il bien comprendre la nature de ces données pour déterminer, en fonction du besoin recherché, le modèle d'association le plus pertinent.

### On distingue quatre grands bénéfices de ces mesures :

### - Une réponse pour accroître l'acuité des mesures sur les plus petites audiences :

Le développement des consommations à la demande favorise la fragmentation des audiences. Au-delà de quelques programmes fédérateurs qui continuent d'attirer massivement le public, l'offre des plateformes à la demande favorise le développement d'une fragmentation croissante des usages. Les panels peuvent continuer à mesurer précisément les audiences les plus importantes, même à la demande, mais la mesure des autres contenus devient un véritable enjeu.

### - Etendre la couverture des usages :

Les mesures d'audiences sont soumises à un ensemble de contraintes de réalisation. Chaque mode de recueil de l'information a ses avantages et ses limites. Ces limites incluent des restrictions du périmètre d'audience mesuré. L'utilisation de data peut permettre d'étendre le champ de mesure d'un système de mesure par panel.

### Des modèles pour répondre aux biais de sélection des panels :

Comme nous l'avons noté précédemment, les panels sont potentiellement soumis à des biais de sélection. L'utilisation de données exogènes exhaustives peut permettre de corriger ces biais par calage sur des niveaux d'usages en complément du calage classique des structures socio-démographiques.

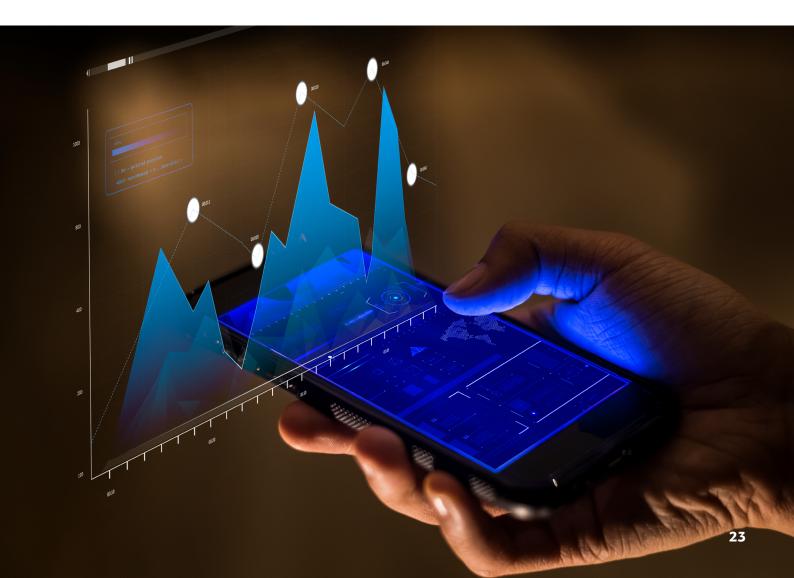
#### - Optimiser les coûts d'une mesure :

La mise en œuvre de systèmes hybrides permet d'envisager une optimisation des coûts des mesures. Sans passer par de tels systèmes, le développement de la précision des mesures par un simple accroissement du nombre d'interviews ou de panélistes aurait un coût beaucoup plus élevé.

Mais une optimisation des coûts ne signifie pas pour autant qu'il sera par exemple possible de baisser le prix d'une mesure existante en remplaçant une partie de l'échantillon par des data. Il y a plusieurs raisons à cela. Tout d'abord, les données des voies de retour ne sont pas gratuites. Les propriétaires de ces données, sauf s'ils sont contraints par voie règlementaire à les livrer gracieusement (ce qui fut par exemple le cas au Canada par suite d'une décision de la CRTC - Canadian Radio-television and Telecommunications Commission), fixent un coût d'accès aux données. Celui-ci peut être significatif. Par ailleurs, les systèmes hybrides fonctionnent par définition sur deux piliers, un issu des données de panels ou d'enquêtes et un autre sur les données voie de retour. Pour être pertinents, les modèles statistiques nécessitent des bases d'apprentissage fiables et conséquentes.

« Pour être pertinents, les modèles statistiques nécessitent des bases d'apprentissage fiables et conséquentes. »

Dans la plupart des cas, il ne s'agit donc pas de baisser des coûts existants, mais d'étendre la finesse et la couverture d'une mesure à un coût bien plus faible que celui qui aurait été nécessaire avec une approche traditionnelle. Les lois statistiques nous apprennent que pour doubler la précision, il faut multiplier par quatre la taille des échantillons et des panels. Dans les faits, cela se traduirait par un accroissement très significatif des coûts de cette mesure. A l'inverse, avec une approche hybride on peut envisager une amélioration de la précision (même bien audelà d'un doublement), pour un coût additionnel moindre.



# Principes & typologie **de mesure hybrides**

### Les grands principes de l'hybridation

La logique de la mesure hybride repose sur quatre grands principes fondamentaux :

### - Complémentarité:

Aucune source ne suffit à elle seule. L'idée est de combler les lacunes d'une source par les qualités d'une autre.

### - Alignement temporel et structurel:

Pour combiner des données, il faut qu'elles soient alignées dans le temps (par exemple, sur la même période) et dans leurs définitions (unités de mesure comparables).

### - Modélisation:

L'hybridation passe par la construction d'un modèle statistique pour combiner les différentes sources de données.

### - Gouvernance des sources :

Une mesure hybride suppose un cadre clair, et en particulier la plus grande transparence sur les données partagées.

> « En Data Science, le terme hybridation n'a pas une définition unique et universelle. »

# Grandes familles de mesure hybride

En Data Science, le terme hybridation n'a pas une définition unique et universelle. Il est utilisé pour désigner des méthodes statistiques de rapprochement de données de sources différentes voire des combinaisons de plusieurs méthodes statistiques. On peut distinguer quelques grandes familles de méthodes qui répondent chacune à des besoins distincts.

### A. Fusion statistique

Cette approche, basée sur les techniques d'imputation, consiste à rapprocher plusieurs bases de données afin de créer un ensemble de données enrichi, cohérent et plus complet. Elle est utilisée lorsqu'aucune source ne contient l'ensemble des variables d'intérêt utiles et que les différentes sources de données ne peuvent pas être directement appariées à l'aide d'un identifiant commun.

L'objectif est de reconstituer une base de données similaire à ce que l'on aurait obtenu si toutes les variables avaient été collectées sur les mêmes individus. Le rapprochement des bases de données s'appuie sur la similarité des individus sur un ensemble de variables communes dans les différentes bases : on définit une distance entre les individus des différentes bases et on les associe ensuite en fonction de leur similarité.

- **Bénéfices:** cette approche permet de limiter le « fardeau de réponse » des panélistes ou interviewés et de reconstituer une base de données complète similaire aux données d'origine et donc facilement exploitable dans des outils de restitution classiques.
- Limites: les variables spécifiques à chaque base de données ne sont jamais observées conjointement. La qualité de la fusion repose donc fortement sur le pouvoir explicatif des variables communes sur les variables spécifiques. En l'absence de variables communes pertinentes, la fusion sera proche de l'aléatoire.

### **B.** Calage

Cette approche est utilisée lorsqu'on dispose d'une source de données issue d'un échantillon ou d'un panel et d'une autre source de mesure exhaustive. Dans ce cas, on veut utiliser l'information issue de la mesure exhaustive, qui correspond à un total connu sur l'ensemble de la population, pour améliorer la précision statistique, réduire la variabilité des résultats ou corriger un biais de sélection sur l'échantillon ou le panel. L'approche consiste à introduire des contraintes de calage supplémentaires dans le redressement de l'échantillon ou du panel.

- **Bénéfices:** cette approche permet d'assurer la cohérence des deux sources de données sans avoir à en modifier la structure. Par ailleurs, elle ne nécessite pas d'avoir accès à la donnée brute de la mesure exhaustive, d'une volumétrie souvent très importante, mais uniquement des totaux sur les variables de calage.
- Limites: les différentes sources de données doivent être parfaitement comparables, ce qui n'est pas toujours nativement le cas. Des prétraitements peuvent donc être nécessaires pour mettre en cohérence les périmètres mesurés et les indicateurs calculés.

### C. Profiling

Cette approche est utilisée lorsqu'on dispose d'une source de données issue d'un échantillon ou d'un panel et d'une autre source de mesure exhaustive et que l'on veut enrichir la mesure exhaustive à l'aide d'informations, souvent très riches, issues d'un échantillon ou d'un panel. En effet, la donnée exhaustive permet d'observer des usages encore rares ou occasionnels qu'un échantillon ne peut mesurer avec précision.

L'approche consiste à construire un modèle statistique de qualification sur les données de l'échantillon ou du panel et de l'appliquer ensuite sur la donnée exhaustive pour l'enrichir.

- **Bénéfices**: cette approche permet d'améliorer la compréhension des usages émergents ou rares sans avoir à augmenter significativement la taille des échantillons.

- Limites: les données exhaustives étant généralement collectées en silo, les variables explicatives à disposition pour la modélisation sont relativement pauvres, ce qui limite la capacité d'un modèle à estimer des profils de manière fiable.

### D. Génération de population synthétique

La génération de population synthétique n'est pas propre aux approches hybrides mais elle peut être utilisée comme étape préalable au rapprochement de différentes sources de données. Elle est particulièrement utile lorsqu'au moins une source de données provient d'une mesure exhaustive. Cette approche est utilisée à l'origine pour l'analyse spatiale à un niveau fin. Elle consiste à construire un ensemble exhaustif et représentatif de la population sur lequel peuvent être distribués des résultats d'enquêtes ou de panels et des data. Cette redistribution pourra faire appel à des méthodes déterministes lorsqu'un identifiant commun est disponible entre les différentes sources ou à des méthodes stochastiques ou probabilistes dans le cas contraire. Les techniques de fusion ou de qualification détaillées précédemment pourront s'appliquer sur une population synthétique, tout comme les techniques de probabilisation.

- Bénéfices: Cette approche permet de combiner différentes sources sans qu'elles soient nécessairement sur des univers strictement comparables et elle facilite la préservation des caractéristiques des données d'origine. Elle permet d'envisager l'exploitation à grande échelle de données individuelles très fines, sans se heurter aux problèmes de gestion de la privacy de ces données.
- Limites: La qualité et la conformité de la population synthétique sont très dépendantes de la quantité et de la granularité des informations mises à disposition par les instituts nationaux de statistique. Ensuite, comme toute approche par modèle, la fiabilité des résultats dépend du pouvoir explicatif des variables communes sur les variables spécifiques.

# Une approche systémique et gouvernée

La mesure hybride n'est pas qu'une affaire de technique : c'est un dispositif systémique, qui suppose des accords avec les acteurs qui partagent leurs données, des choix méthodologiques communs, des standards transparents. Elle exige une gouvernance forte :

- Interopérabilité des formats et des définitions.
- Transparence et auditabilité des modèles par des tiers indépendants.

C'est ce qui distingue une vraie mesure hybride d'une simple juxtaposition de chiffres.

La mesure hybride n'est ni une fin ni une mode: c'est une nécessité dans un monde média devenu trop fragmenté pour être saisi par un seul instrument. Elle permet d'articuler des visions différentes, de construire une mesure plus complète, plus nuancée, plus crédible.

Encore faut-il qu'elle soit intelligible. Le défi des années à venir est autant technologique que pédagogique: expliquer les modèles, documenter les méthodes, rendre lisible une métrologie de la complexité.

# Les défis inhérents aux mesures hybrides

Les bénéfices des mesures hybrides sont incontestables et nombreux comme nous venons de le voir. Plus grande granularité, stabilité des résultats, couverture de la mesure, optimisation des coûts.... sont autant de bénéfices indéniables qui deviennent indispensables compte tenu de l'évolution des usages médias. Il faut néanmoins noter que les modèles & approches hybrides sont toutefois confrontés à un ensemble de défis.

#### Accès aux données

Ce point est probablement un des plus grands défis de la construction des mesures hybrides. Les données sont partout, mais leur recueil est complexe. Il est soumis à des contraintes juridiques très strictes notamment dans un cadre nécessitant le rapprochement de données issues de différentes sources. Ces données peuvent avoir une grande valeur pour les acteurs qui la collectent. Au-delà d'une question de coût d'accès de la donnée, de nombreux acteurs refusent le partage de leur donnée. Ces données peuvent aussi représenter un enjeu pour les éditeurs et nécessiter des accords contractuels complémentaires entre l'acteur mesuré, l'acteur qui collecte et le mesureur. La législation évolue progressivement sur ce sujet et un nombre croissant de textes européens font référence à la nécessité, sous certaines conditions, du partage de données avec un tiers mesureur, neutre et indépendant. C'est notamment le cas du texte European Media Freedom Act. A ces enjeux juridiques et contractuels s'ajoutent potentielles complexités d'accès données selon les configurations techniques et d'architectures propres à chaque acteur. Un autre défi est la pérennité de l'accès aux données, soumis à la fois aux évolutions des contraintes juridiques, des stratégies des différents acteurs et des technologies.

#### Cohérence des données

(en termes de périmètre couvert ou de définition des indicateurs) Les approches hybrides amènent à rapprocher différentes sources de données panel et voie de retour, mais aussi différentes sources de données voie de retour entre-elles. La mise en œuvre de ces combinaisons de données peut perdre tout son sens si un travail méticuleux d'ajustement des périmètres, et définitions, n'est pas réalisé en amont. Vouloir calibrer un panel de mesure de la télévision à domicile (incluant des postes de TV connectés et des postes non connectés) avec une donnée voie de retour dont les résultats intègrent les usages de chaînes repris uniquement dans l'environnement de la CTV est voué à un échec certain. Le modèle hybride produira bien des résultats, mais qui seront sérieusement biaisés.

### Normalisation des données

Une fois le problème de l'accès aux données résolu, le problème de la normalisation et certification des données se pose. C'est un enjeu important d'autant plus qu'il concerne des acteurs locaux et internationaux. De nombreuses questions se posent : Qui définit les normes et le cadre de l'audit ? Qui le finance ? Qui le réalise ? L'existence et la qualité de telles normes et de tels audits sont essentielles pour pouvoir produire une mesure hybride.

#### Taille des échantillons

Par nature, les systèmes hybrides ont notamment été construits pour dépasser les limites des échantillons en termes de granularité. Dans l'univers des mesures hybrides, on peut résumer en disant que les données census fournissent les niveaux de consommation, les données de panel, permettent d'apporter la dimension individuelle, le profil, et la duplication. La qualité des données panel est donc un élément fondamental de toute mesure hybride.

#### Connexion des données

certain nombre d'approches hybrides nécessitent de connecter les Data d'éditeurs avec celle des panels. Dans ce contexte d'hybridation, la mise en place de data clean room (environnements sécurisés permettant le croisement de données) semble être une réponse adaptée pour enrichir données panels datas disponibles. Cette technique nécessite évidemment le recueil préalable du consentement auprès des utilisateurs (panel et data). Ce travail représente un effort significatif (sur le plan technique et juridique), et ce d'autant que le nombre de sources de données à agréger est grand.

### Richesse des données

Dans les méthodes probabilistes d'hybridation, il faut pouvoir s'appuyer sur des variables communes pour faire le rapprochement entre les data. La quantité de variables communes et leur pouvoir explicatif auront une influence sur la qualité du modèle hybride.

### Intelligibilité des modèles et cohérence des résultats

Les mesures hybrides font intervenir des modèles statistiques qui reposent sur des hypothèses, parfois difficilement vérifiables. De plus, même si l'opérateur de la mesure fait preuve d'une grande transparence sur les méthodes utilisées, certains types de modèles, notamment les modèles de Deep Learning, sont difficiles à interpréter car on ne peut exprimer simplement les relations entre les variables explicatives et celles à expliquer. Par ailleurs, contrairement à la donnée panel qui repose sur l'observation de comportements réels, tout modèle est une vision simplifiée de

la réalité qui ne peut donc pas capturer toute la diversité et la complexité de la réalité. L'adoption d'une mesure hybride par un marché repose donc sur la confiance de l'ensemble des acteurs.

### Estimation de la précision

On peut estimer de manière assez simple la précision statistique des résultats d'enquêtes et calculer des intervalles de confiance. Mais lorsqu'il s'agit de résultats issus d'une mesure hybride, l'estimation est complexe voire impossible. L'analyse de l'évolution des résultats dans le temps est donc plus délicate et nécessite de mobiliser des outils différents des tests statistiques classiques.

### Echantillonner la data

Associer les deux mots «data» et «échantillonnage» peut sembler paradoxal, mais cela reflète une nouvelle réalité des outils mesure «analytics».

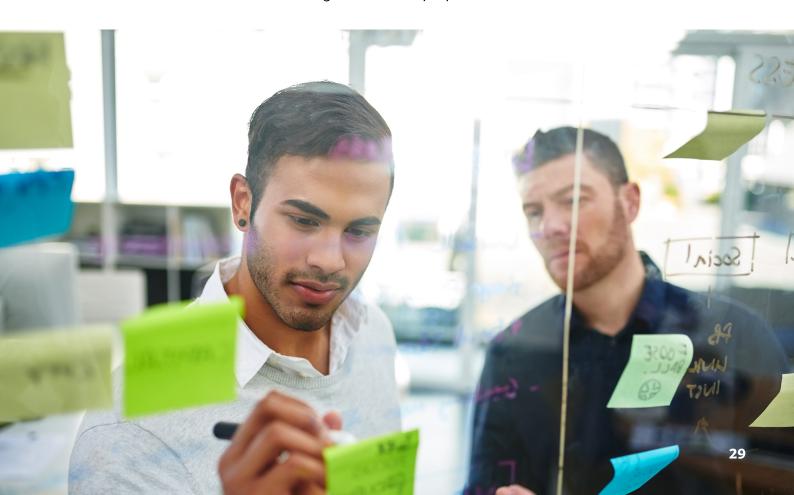
Ces systèmes de mesure ont été conçus pour fournir des rapports et des analyses d'usage avec le plus haut niveau de granularité, puisqu'ils se proposent de capturer et de rapporter de manière exhaustive tous les événements d'usage.

Ces systèmes, comme nous l'avons expliqué, répondaient au besoin d'une nouvelle génération de mesures à l'ère du ciblage, de la personnalisation et de la fragmentation de l'audience.

Mais la collecte de millions, voire de milliards d'événements d'usage a un coût financier et un coût carbone importants : les coûts de stockage et de traitement de ces données. Même avec la baisse des coûts de traitement grâce au cloud, l'utilisation de toutes les informations disponibles pour calculer certains rapports complexes peut s'avérer onéreuse. De plus, pour les outils de reporting, le traitement de l'exhaustivité des données peut nécessiter plus de temps de génération.

« La collecte de millions, voire de milliards d'événements d'usage a un coût financier et un coût carbone importants »

Enconséquence, de nombreux outils d'analyse ont mis en œuvre des algorithmes d'échantillonnage. Dans de nombreux cas, les rapports fournis à l'utilisateur ne sont pas basés sur la totalité des événements, mais sur un échantillon de ceux-ci. Les outils affichent généralement le taux d'échantillonnage et la marge d'erreur que cela implique.



# Les exemples de mesure hybride dans le monde

De nombreux systèmes hybrides alliant données et panels ont vu le jour à travers le monde. Les exemples suivants offrent un aperçu non exhaustif de ces initiatives.

# L'hybride en France dans les mesures de Médiamétrie

A date, des systèmes hybrides combinant data et panels ont été développés pour un certain nombre de mesures d'audience.

### **Mesure Internet (Internet Global)**

Depuis une dizaine d'années, la mesure Internet Global de Médiamétrie NetRatings utilise une base de données du volume de visites de plus d'une centaine de sites et applications pour calibrer les résultats d'audience issus de son panel. Les données en entrée de l'hybridation issues d'outils analytics tiers (ex: Piano, Wysistat, NSP, Experian,...) sont certifiées mensuellement par l'ACPM (Alliance des Chiffres pour la Presse et les Médias). Le périmètre hybridé concerne tous les écrans digitaux, tous lieux, en France métropolitaine. La publication des audiences hybridées est mensuelle, elle concerne plus de 5000 marques et 1000 applications.

### **Mesure Cross-Media Publicitaire (Watch > pub)**

publicitaire Cross-Vidéo mesure Médiamétrie intègre depuis janvier 2025 un système de production hybride qui repose sur les volumes d'impressions issus des données des AdServers afin de mesurer l'apport de la CTV aux audiences linéaires. Pour accélérer l'évolution de solution publicitaire Cross-Média Vidéo, Médiamétrie a annoncé en septembre 2025 un partenariat avec Audience Project. L'accord va permettre de combiner l'expertise méthodologique et les données propriétaires de Médiamétrie avec les actifs technologiques d'Audience Project, déjà intégrés aux principales plateformes du marché. Prévue pour le 1er trimestre 2026, la mesure publicitaire Cross-Média Vidéo proposera une évaluation « précise » de la couverture dédupliquée et de la répétition des campagnes sur tous les canaux, y compris les grandes plateformes, les réseaux sociaux, la télévision connectée et la vidéo en ligne ; et une vision unifiée des performances des campagnes.

### TV segmentée / Adswitching

Dans le cadre du mediaplanning de la télévision linéaire, un système a été développé pour permettre une correction hebdomadaire des GRP sur écran TV, en soustrayant les impressions diffusées via la télévision segmentée (Adswitching). Ce système utilise les logs et les impressions des serveurs publicitaires (AdServers) pour effectuer cette correction.

# Les autres projets de mesure hybrides de Médiamétrie :

Dans le cadre d'autres produits et services, Médiamétrie peut également s'appuyer sur des données externes afin de compléter, enrichir et optimiser les dispositifs de mesure existants. Ces données permettent d'élargir les périmètres couverts, d'affiner les indicateurs, et de répondre à des besoins spécifiques du marché en matière de précision, de granularité ou de fréquence.

POC d'hybridation de la mesure des chaînes thématiques à partir des données de boxes (Proof of Concept de la mesure TV) Médiamétrie a mené entre 2021 et 2024 un proof of concept (POC) en collaboration avec les principaux opérateurs français - Orange, SFR, Bouygues Telecom et Canal+ - afin d'explorer l'intégration des données de retour des boxes dans la mesure

du live des chaînes thématiques. Pour cela, Médiamétrie a recruté un panel de foyers au sein de chaque opérateur et mis en place une collecte quotidienne des data opérateurs. Un important travail a été réalisé pour uniformiser les data opérateurs via le traitement de leurs données (nettoyage, écrêtage...). En parallèle, des chantiers méthodologiques ont été menés par Médiamétrie, incluant des itérations de modèles, la multiplication du panel Médiamat, des opérations de fusion et d'individualisation.

Les données hybrides produites ont démontré de nombreux avantages. Les niveaux mesurés se sont avérés très proches de ceux de la mesure de référence TV. De plus, le nombre de cases à audience nulle ainsi que la stabilité globale ont été nettement améliorés. Mais ce POC a aussi permis de mettre en lumière les enjeux majeurs liés à une mesure hybride, notamment :

- Les aspects juridiques, avec des exigences RGPD et la nécessité d'intégrer de multiples partenaires,
- La complexité à normaliser la qualité des données en entrée.
- La complexité du paysage audiovisuel français, marquée par une multiplicité d'acteurs, de modes de consommation et des offres qui se recouvrent, ainsi que par une part non négligeable de réception TNT sans voie de retour,
- -La capacité du marché à accepter les impacts d'une rupture méthodologique,
- La question du financement, en particulier lorsque les données externes doivent être acquises ou valorisées.

En octobre 2024, Médiamétrie, en concertation avec les acteurs du marché, a dressé le bilan du POC et a pris la décision de ne pas poursuivre les explorations engagées sur l'univers des chaînes thématiques.

De nouveaux travaux exploratoires sont en cours pour poursuivre l'objectif d'hybrider les résultats de la mesure d'audience de la télévision. Le POC se limitait aux données de la TV linéaire. D'autres sources de données et méthodes sont actuellement à l'étude afin de bâtir un système hybride pour la mesure des audiences à la demande.

### Projet de mesure Radio digitale hybride (EAR Insights Digital)

Des tests sont en cours pour produire une mesure qui combine le panel de mesure des audiences EAR Insights utilisé pour le médiaplanning radio avec les données certifiées mensuellement par l'ACPM sur la consommation du live digital des radios et webradios (données issues des Content Delivery Network utilisés pour diffuser l'audio live en streaming).

### Mesure hybride des plateformes (Watch)

La mesure Watch produite par Médiamétrie repose sur des données de panel collectées à domicile via un audimètre de streaming. Cette mesure collecte les données via l'usage du Wifi à domicile sur tous les écrans connectés pour les audiences des plateformes de streaming.

Pour étendre le périmètre à tous les lieux et à tous les écrans, Médiamétrie proposera, après une période de test, une approche hybride incluant les logs et le volume de consommation globale au second semestre 2025. Pour réaliser cette hybridation, deux sources de données peuvent être utilisées: soit les mesures issues du SDK, soit les données «server side» fournies par les plateformes (ce qui nécessite un audit pour valider le périmètre mesuré).

# L'hybride à l'international : quelques exemples clés

# Les exemples relatifs à la mesure TV

### **Autriche - AGTT**

L'Autriche a lancé en septembre 2024 une nouvelle mesure d'audience de la télévision qui fonctionne grâce à une combinaison de données de panels traditionnels et de data de TV connectées via la technologie HBBTV. Les données digitales proviennent de la collecte de données de 100 000 participants à un panel synthétique et de 1 000 000 téléviseurs connectés et décodeurs. A ces données s'ajoutent celles d'un panel de 3 200 participants à une mesure TV. Ces dernières informations servent à calibrer et vérifier les modèles. Cette approche hybride permet de réduire les fluctuations de résultats, mais aussi de réduire la part des résultats qui, par manque de finesse de la mesure, se traduisent par l'affichage d'un zéro. Ce système d'audience est couplé aux AdServers (serveurs publicitaires) des chaînes pour permettre un ciblage publicitaire en temps réel (sur le modèle de la publicité segmentée en France).

Cette nouvelle mesure constitue la nouvelle référence officielle en termes d'audience.

#### Italie - Auditel

Auditel a commencé à mesurer l'audience de la télévision en 1984. La taille de l'échantillon a été doublée en août 1997 et triplée en juin 2017 pour répondre aux enjeux de la fragmentation croissante des audiences : le SuperPanel™ est aujourd'hui composé de 16 100 foyers.

Trois objectifs principaux ont été visés avec ce projet :

- **1.** une plus grande granularité dans la mesure des petites chaînes thématiques et locales ;
- **2.** une plus grande stabilité dans la mesure des chaînes plus importantes : la taille accrue de l'échantillon a permis de réduire les facteurs d'extrapolation des deux tiers ainsi que l'intervalle de confiance ;
- **3.** Améliorer la représentation fidèle de la diversité de la société italienne.

L'échantillon qui constitue l'Auditel SuperPanel™ est la somme de deux échantillons : un panel individuel (People Meter Panel - PM) et un panel foyer (Set Meter Panel - SM). Le People Meter Panel (PM) compte 5 682 foyers. Le Set Meter Panel (SM) compte 10 294 foyers proportionnels. Le SM comme le PM incluent des foyers avec et sans téléviseur. Les données du Set Meter Panel sont collectées sans préciser qui regarde la télévision (il n'y a pas de télécommande pour déclarer sa présence devant l'écran). Cette information est modélisée à l'aide des données du People Meter Panel.

### Suisse - Mediapulse

En Suisse, Mediapulse a lancé le 1er juillet 2022 un système hybride de mesure de l'audience télévisuelle pour répondre notamment aux attentes d'une plus grande finesse des résultats.

Le système de Mediapulse combine deux sources principales de données :

### - Panel de mesure TV:

Un échantillon représentatif de ménages suisses (1870 ménages), équipé de dispositifs de mesure, fournit des données détaillées sur la consommation télévisuelle.

### - Données des décodeurs numériques :

Les données d'environ160 000 décodeurs fournis par des opérateurs tels que Swisscom et Sunrise (ex UPC) transmettent des informations sur l'utilisation de la télévision. Ces données sont anonymisées et modélisées pour créer un «panel virtuel» de 15 000 profils.

Ces deux sources de données sont combinées par modélisation pour créer une mesure plus fine. Le système hybride concerne les données de la TV linéaire ainsi que la publicité digitale sur CTV.

Les résultats de la mesure sont d'abord publiés à J+1, en reposant uniquement sur les résultats du panel. Lors du lancement du système, la version hybridée était publiée à J+7, en même temps que les résultats consolidés du Replay. Elles sont désormais publiées à J+3. Notons que ces données constituent la référence. Les données publiées le lendemain de la diffusion sont dites « temporaires ».

Les données hybrides sont utilisées comme référence pour le planning, le reporting et les transactions depuis 2022.

### **Espagne**

En Espagne, Kantar Media est en train de passer à un nouveau système de mesure hybride combinant un panel avec big data. Développé en collaboration avec Konodrac, fournisseur de mesures HbbTV, ce nouveau système fusionne les données census de la Télévision linéaire issues des téléviseurs connectés avec les données du panel afin d'améliorer la stabilité des résultats et de réduire les audiences nulles.

La mesure est auditée par l'autorité AIMC en Espagne. Cette commission de contrôle qui régit la mesure de référence de la télévision en Espagne fournie par Kantar Media est actuellement en train de finaliser son audit du passage à ce nouveau système.

### Canada - Numeris

Au Canada, le gouvernement, dans une décision de la CRTC, a imposé aux câblo-opérateurs de mettre à disposition de Numéris, le mesureur local de la télévision, des données issues de leurs boxes. Numeris travaille depuis quelques années à la production d'une mesure complémentaire à la mesure TV de référence. Cette mesure est à ce stade publiée en parallèle de la mesure de référence.

En avril 2024, Numeris a lancé sa nouvelle mesure hybride de la télévision linéaire, baptisée

Enhanced Television Audience Measurement (Enhanced TAM ou eTAM). Cette approche combine les données traditionnelles des panels avec des données provenant des boîtiers des décodeurs (set-top boxes).

La mesure Enhanced TAM se base sur les sources suivantes :

#### - Panel de mesure TV :

Basé sur 4500 foyers, 11 000 individus panélistes via le PPM (Portable People Meter).

### Données voies de retour des set top boxes (opérateurs) :

Ces données proviennent d'échantillons de dizaines milliers de foyers sélectionnés parmi les abonnés des différents opérateurs. Il ne s'agit donc pas d'un ensemble exhaustif des données voies de retour.

Ces deux sources de données sont ensuite combinées par fusion. L'échantillon hybride est limité en taille du fait des limites des outils de restitution.

Les bénéfices de cette approche hybride sont les suivants :

- Une mesure d'audience plus complète, grâce à l'intégration des données des chaînes disponibles via les câblo-opérateurs, qui ne sont pas nécessairement couvertes par le panel traditionnel.
- Une plus grande stabilité des résultats, rendue possible par un volume d'observations plus important,
- La baisse ou disparition des audiences à zéro pour certaines chaînes, tranches horaires ou cibles,
- Des gains en termes de coût via la modélisation de panélistes sans avoir à les recruter directement dans le panel.

Les résultats d'audience issus de cette hybridation sont disponibles dans une base distincte à visée analytics. A date Enhanced TAM n'est pas la donnée de référence, les données sont utilisées en parallèle du TAM classique, l'adoption de l'hybride en tant que currency de référence est en cours avec le marché.

Numeris a également développé une mesure d'audience Vidéo (VAM). Ce service couvre la consommation de contenu vidéo sur la télévision linéaire, les services OTT, la SVOD, le streaming...

En septembre 2025, Numeris va déployer cette solution à l'échelle nationale, offrant ainsi une vue d'ensemble plus complète des habitudes de consommation vidéo à travers le pays. VAM repose également sur un système hybride (Panel + data) en partie imbriqué avec Enhanced TAM. En effet, Enhanced TAM permet de faciliter une meilleure fusion et attribution entre les audiences de la télévision linéaire et de la vidéo numérique dans le cadre de VAM (Video Audience Measurement).

### Etats-Unis - Nielsen

Nielsen a publié récemment une nouvelle mesure de la télévision : Big Data + Panel. Ce dispositif a été audité de manière complète par le MRC en janvier 2025, et les données sont utilisées comme mesure d'audience pour la saison 2025/2026, qui commence en septembre 2025. La mesure traditionnelle, basée uniquement sur un panel de foyers et d'individus, est disponible pour la mesure d'audience jusqu'à la fin de 2025.

La mesure Big Data + Panel de Nielsen se base sur les sources suivantes :

Panel de mesure TV: 42 000 foyers, soit plus de 101 000 individus.

#### Sources pour le Big Data:

- Données anonymisées provenant de millions de décodeurs TV, téléviseurs connectés et plateformes de streaming,
- Données voie de retour (RPD) des décodeurs câbles et satellites,
- Données de reconnaissance automatique des contenus (ACR) des Smart TV,
- Données propriétaires de services de streaming participants pour évènements en live (ex : Amazon Prime pour la NFL).

Le modèle d'hybridation est basé sur une attribution de l'audience device aux individus et une projection sur l'ensemble des télés pectateurs.

L'ensemble des indicateurs délivrés dans le cadre de la mesure par panel seul sont disponibles dans la mesure Big Data+ Panel. Les données hybridées sont toutefois livrées un jour plus tard (à J+2) que les audiences du panel seul. La mesure hybride aux États-Unis concerne principalement l'écran de TV, l'intégration des autres devices est en développement mais pas encore généralisée

#### Suède - MMS

La Suède a été l'un des premiers pays à exploiter des données voie de retour pour développer la mesure de la télévision de rattrapage. En 2011, MMS a lancé la première mesure de la BVOD en utilisant des données analytics des éditeurs. Au fil des ans, ils ont développé de nouveaux services et de nouvelles capacités et ont connecté ces données avec celles de leur panel TV.

### Afrique francophone - Médiamétrie / Canal+

d'Afrique francophone, pays Médiamétrie a déployé pour le compte de son client Canal+ une mesure d'audience innovante fondée sur une approche hybride. Le dispositif repose sur la collecte des logs techniques issus des boxes que possèdent un panel de plus de 8000 foyers abonnés Canal+ recrutés dans les principales villes des 9 pays couverts par l'étude, enrichis par une enquête déclarative menée auprès d'une partie de ces foyers. Les boxes des foyers recrutés sont équipées par les équipes terrain de clés USB 3G/4G capables de transmettre des données d'usage (changement de chaînes, horaires, etc.). Ces données d'usage sont ensuite transformées en audience téléviseur, puis en audience individuelle grâce à un modèle statistique avancé, basé sur les chaînes de Markov cachées. Ce modèle tient compte de la composition des foyers (notamment les liens de parenté entre chaque membre, du nombre de personnes, etc.), des habitudes de visionnage déclarées par chacun des membres, et du format des différentes chaînes qui composent l'offre de télévision dans les pays couverts.

Cela permet d'estimer avec précision qui regarde quoi, et quand. La mesure couvre plus de 240 chaînes, avec une granularité à la seconde, et fournit chaque semaine des résultats en jour daté, au niveau national et aussi panafricain. Ce dispositif permet à Canal+ de disposer, d'une vision fine des comportements d'audience,

essentielle pour piloter ses contenus, ses stratégies de diffusion et la commercialisation de son espace publicitaire. C'est une avancée majeure sur ce territoire avec des applications concrètes inédites telles que des analyses d'audience sur un jour et un programme en particulier, des plans médias optimisés au jour près, ou encore des bilans de campagnes après la diffusion des campagnes à destination des agences et annonceurs.

### Autres projets en cours, notamment en Europe

Le sujet de la mesure hybride TV n'est pas limité à ces quelques pays, le sujet est en cours de discussion dans de nombreux pays ou des tests / POCs sont menés.

### Allemagne - AGF

En Allemagne, l'AGF a annoncé dès 2024 un projet avec l'acteur UTIQ pour créer un système qui permet de marier des données de panel avec un graph d'identifiants à grande échelle (UTIQ a développé une solution en utilisant les données de plusieurs grands opérateurs télécoms en Europe).

### Royaume-Uni - BARB

Le BARB qui opère la mesure de la télévision au Royaume Uni a lancé un appel d'offres en mai 2024 pour <u>développer l'intégration de donnée voie de retour dans sa mesure de la télévision</u>. Le projet BARB Panel Plus comporte trois volets :

### - Méthodologie:

Conception d'une méthode d'hybridation adaptée afin de réunir les données de panel avec les données de voie de retour, les données propriétaires de serveur et les données de recensement BVOD.

- Traitement et livraison des données: Développement d'un logiciel permettant de traiter les données selon la méthodologie d'hybridation, selon un calendrier et à la demande.

#### - Fourniture de données :

Sources de données de consommation TV, telles que les données de voie de retour et les données propriétaires de serveur, à hybrider les données du panel.

Par ailleurs, le BARB propose depuis plusieurs années un système d'hybridation des données de la BVOD (replay) : c'est le projet Dovetail. Le modèle Dovetail repose sur les données census de la BVOD, tandis que le projet Barb Panel Plus vise à enrichir cette approche en intégrant également les données de voie de retour issues des boxes et Smart TV, ainsi que des données propriétaires (first-party) des éditeurs.

### Le cas de la mesure Cross-Médias

### WFA - Projets Origin (UK) / Aquila (Etats-Unis)

En 2019, la WFA a édité un manifeste sur la mesure Cross-Médias. Celui-ci définit un ensemble de principes clés repris depuis par une grande partie de la profession. Ce manifeste donna par ailleurs lieu au lancement d'un grand projet (décliné dans une version au UK, Origin, et aux Etats-Unis, Aquila) qui met en œuvre les grands principes de cette nouvelle mesure.

Au cœur du système de mesure, une approche nouvelle de traitement de la donnée d'audience, le Virtual ID model, qui résulte de l'intégration et de la comparaison des données provenant des serveurs publicitaires de différents partenaires et d'un panel de haute qualité, sans générer de problèmes liés à la privacy. C'est ce que l'on appelle les modèles de population synthétique. Jusque-là, lorsqu'une donnée analytique était utilisée dans un projet de mesure hybride, on cherchait à retirer tous les éléments sensibles de cette donnée qui permettrait l'identification des individus (comme l'adresse IP par exemple, ou des éléments de profils précis) dont elle est issue.

A l'inverse dans un modèle de population synthétique comme celui-ci, on ne détruit pas ces informations, on les transforme et on les assigne à des individus virtuels, miroirs de chaque individu de la population mesurée. Ainsi

dans le projet Origin, il existe autant de virtual ID qu'il y a de personnes vivant au Royaume Uni. Cette transformation garantit l'anonymat des données et permet de respecter les législations de protection de la privacy tout en conservant la finesse de la donnée en entrée. Chaque log de données permet ainsi d'alimenter une population virtuelle dont l'analyse ouvre la voie à de nouvelles générations de mesure.

La mise en œuvre des projets Aquila & Origin ont nécessité1 des investissements importants pour concevoir les systèmes techniques qui permettent la collecte et le traitement de ces données.

Compte tenu de son ambition et de son envergure, le projet est souvent considéré comme la « North Star » en matière de mesure Cross-Médias.

### Finlande - SpotOn

SpotOn est une devise publicitaire vidéo multiplateforme en Finlande qui combine la télévision linéaire, le streaming et la vidéo en ligne. Sa méthodologie s'articule autour d'un panel de référence à source unique créé en enrichissant le panel TV national avec les impressions du serveur publicitaire liées à ses foyers, ses membres et ses appareils. Ce panel à source unique est ensuite utilisé comme « terrain d'apprentissage » pour modéliser chaque étape du processus de mesure.

Le système est conçu pour répondre en temps quasi réel aux demandes des utilisateurs en adaptant une cascade de modèles à chaque tâche.

Les résultats sont fournis via une application web ou une API et comprennent des estimations précises des contacts, de la portée et de la fréquence dédupliquées, ainsi que de l'accumulation de la portée pour tout groupe cible d'âge et de sexe requis.

Le système de mesure SpotOn est détenu et produit par JIC Media Metrics Finland Oy (MMF).

### International - Kantar Media Campaign Audience Validation (« CAV »)

En 2021, Kantar Media a lancé Campaign Audience Validation (« CAV ») au Brésil et en Colombie, puis dans cinq autres marchés en Amérique latine, en Europe et en Asie du Sud-Est, avec d'autres marchés en préparation. Ce service fournit quotidiennement aux annonceurs des indicateurs de performance d'audience des campagnes, tels que la portée et les données démographiques, actuellement pour la télévision linéaire et le digital. Pour ce faire, il utilise les données d'exposition publicitaire linéaire provenant du linear currency service et reçoit les données d'exposition publicitaire des plateformes numériques par différentes méthodes, notamment l'import direct et les échanges de panels. La solution est adaptée à chaque marché local, en fonction des exigences et des ressources disponibles sur place.

### Les exemples relatifs à la mesure Radio

#### Australie - GfK Radio360

En Australie, une nouvelle mesure de référence de l'audience de la Radio a été lancée en juin 2023. Cette mesure, Radio360, repose sur un dispositif hybride combinant trois sources de données:

- Panel de mesure déclarative par carnets d'écoute hebdomadaires en ligne, auprès de 50 000 personnes chaque année, pour recueillir l'audience globale de la radio et des stations,
- Panel de mesure automatique de 2 000 panélistes équipés de la montre GfK MediaWatch,
  Des Tags Gfk et logs des radios pour mesurer l'écoute en direct sur les supports numériques (streaming).

Le pilier de cette mesure reste le carnet d'écoute, qui fournit le niveau d'audience global des stations et du média. Les données voie de retour et les données issues de la mesure automatique MediaWatch permettent d'affiner les résultats selon les modes d'écoute.

Les données voie de retour, fournies directement

par les éditeurs (logs) ou collectées via un tag GfK, sont profilées grâce à la mesure MediaWatch. Les données ainsi hybridées déterminent la contribution et le profil de l'écoute online (streaming) au sein du niveau global des radios mesuré par carnets d'écoute. Les niveaux d'audience globaux issus du carnet d'écoute sont ainsi préservés.

# Les exemples relatifs à la mesure du Digital

# UK - Ipsos iris / UKOM

Au Royaume-Uni, la mesure digitale du UKOM est produite par Ipsos. Le système repose sur trois grandes sources de données: une étude de cadrage, un panel de 10 000 individus avec un système de mesure passive et des données site-centric des sites participants intégrant les données d'environ 100 millions de terminaux. La mesure site-centric est basée sur des cookies tiers associés à une solution d'empreinte digitale afin de faciliter l'identification transversale des devices. Les devices des panélistes sont également identifiés permettant ainsi de faire le lien entre données panel et données site-centric. La méthodologie consiste à créer à partir de ces sources de données, une population synthétique d'environ 1 million de personnes dont les usages représentent ceux de la population du Royaume-Uni.

## International - Nielsen DAR (Digital Ad Ratings)

Ce service mesure l'audience des campagnes numériques à l'aide des tags Nielsen et d'une intégration côté serveur avec les données des plateformes et de partenaires DATA, ainsi que des mesures démographiques basées sur des enquêtes et des panels afin d'améliorer la précision des données démographiques attribuées. Cette mesure est utilisée pour communiquer l'audience (couverture, données démographiques) des campagnes digitales. L'association de multiples sources de données au sein de cette mesure en fait une illustration claire des méthodologies hybrides.

---

# L'intelligence artificielle: de nouvelles perspectives pour les mesures

# Une **nouvelle ère** pour la mesure d'audience ?

# « L'IA est souvent évoquée comme la promesse d'un champ des possibles élargi »

L'intelligence artificielle (IA) bouleverse aujourd'hui l'ensemble des secteurs économiques, et celui de la mesure d'audience n'y échappe pas. Dans un environnement média de plus en plus fragmenté, où les données sont massives, hétérogènes, et souvent incomplètes, l'IA est souvent évoquée comme la promesse d'un champ des possibles élargi et de perspectives d'analyses nouvelles. Ce chapitre explore comment l'IA pourrait transformer la mesure d'audience, en s'appuyant sur des cas d'usage concrets et en posant les questions méthodologiques et éthiques que cela implique.

# IA: de quoi parle-t-on?

Il existe de multiples définitions de l'IA, qui dépendent du domaine d'expertise, des objectifs et des usages associés.

On présente parfois l'IA comme une version augmentée de la Data Science. L'IA et la Data Science sont en réalité deux disciplines distinctes même si elles reposent sur des fondements théoriques communs et ont des zones de recouvrement, comme le machine learning par exemple. Elles se distinguent principalement par leur finalité: la Data Science a pour objectif de comprendre, d'expliquer, d'analyser ou de prédire alors que la finalité première de l'IA est de reproduire des tâches complexes à l'intelligence humaine. Ces disciplines sont donc complémentaires et l'IA n'a pas vocation à remplacer les méthodes statistiques classiques, elle démultiplie les capacités d'analyse par l'automatisation des tâches complexes, facilite le traitement de volumes massifs de données

et ouvre la voie à des modèles prédictifs plus performants.

Règlement européen l'intelligence sur artificielle (IA Act) définit un système d'intelligence artificielle (SIA) par sa capacité d'inférence, c'està-dire sa capacité à générer des prédictions, du contenu, des recommandations ou des décisions, et sa capacité à inférer des modèles ou des algorithmes, ou les deux, à partir de données. Les techniques permettant l'inférence lors de la construction d'un système d'IA comprennent des approches d'apprentissage automatique et des approches fondées sur la logique et les connaissances qui font des inférences à partir des connaissances encodées ou de la représentation symbolique de la tâche à résoudre.

En Computer Science, l'IA désigne un ensemble de techniques et algorithmes visant à automatiser des tâches associées à l'intelligence humaine, comme le raisonnement, l'apprentissage ou la compréhensiondulangage. On peut citer quelquesunes des principales familles de techniques et leurs objectifs et domaines d'application.

# 1. Apprentissage automatique (Machine Learning)

L'apprentissage automatique est un domaine de l'intelligence artificielle qui consiste à apprendre à partirde données sans avoir à expliciter de règles ou relations a priori. C'est ce qui distingue les modèles de machine learning des modèles statistiques classiques. Selon la nature des données, différentes techniques peuvent être utilisées:

L'apprentissage supervisé lorsque les données d'apprentissage sont labelisées ou étiquetées, par exemple pour qualifier et enrichir de données logs clients à partir de données panels. L'apprentissage non supervisé lorsque les données d'apprentissage ne sont pas étiquetées et que l'on cherche à découvrir des relations cachées, par exemple segmenter des comportements, détecter des anomalies ou identifier des motifs récurrents dans les données.

# 2. Apprentissage profond (Deep Learning)

L'apprentissage profond est une sous-catégorie de l'apprentissage automatique basée les réseaux de neurones, modèles inspirés du fonctionnement du cerveau humain, composés de couches de neurones artificiels qui transforment et apprennent à partir d'un ensemble de données. Ce mécanisme permet de traiter de gros volumes de données complexes ou non structurées (images, sons, textes), notamment dans les analyses de contenus, et de résoudre des tâches comme la traduction ou la reconnaissance d'images par exemple. L'apprentissage profond utilise des réseaux de neurones comportant nombreuses couches cachées utilisé par exemple pour la classification ou modélisation de fonctions complexes.

# 3. IA générative

L'IA générative est un domaine particulier de l'IA consacré à la génération de contenus (texte, image, son, vidéo...). Elle s'appuie sur le deep learning et le natural language processing (NLP) pour apprendre sur des corpus de données massifs et créer du contenu nouveau. Les modèles d'IA générative existent depuis plusieurs années. Leur utilisation a cru de façon exponentielle depuis la découverte par le grand public de ChatGPT 4. Le nombre d'applications et d'utilisations de ces outils s'est depuis démultiplié (création de code, d'images, de vidéos, de musique, de présentations...).

# IA dans les études

Les outils d'I Apermettent d'automatiser des tâches répétitives des différentes étapes de réalisation d'une étude. Ils peuvent faciliter notamment la préparation des questionnaires (création du questionnaire/guide d'entretien, traduction), la collecte des données (chatbot, agent IA...), l'analyse des données (codification, transcription, traitement, modélisation) ou la restitution.

La déclinaison de chabots développés pour collecter des données pose encore de nombreuses questions et se heurte à de nombreuses limites. Si la génération vocale et la qualité d'expression atteint désormais un très haut niveau, il reste impossible d'envisager de confier à des machines des questionnaires CATI, que ce soit en raison de l'importance de l'accroche initiale pour convaincre un interlocuteur de répondre ou de celle de la complexité d'un questionnaire d'audience.

Sur le volet de la restitution, l'IA générative ouvre en revanche des perspectives prometteuses. On peut désormais envisager d'accompagner chaque livraison de résultats de synthèses automatiques qui en extraient les points clés: les principales évolutions mises en perspective, les tendances marquantes identifiées, les éléments nécessitant une attention particulière signalés. Cette première lecture augmentée permet ainsi d'améliorer l'efficacité du processus d'analyse des résultats. L'IA générative pourra s'intégrer également dans les outils de restitution sous la forme de chatbots pour permettre aux utilisateurs d'interroger l'interface par des questions formulées en langage naturel plutôt que par des requêtes parfois complexes.

En 2024, Ipsos a ainsi lancé un outil d'analyse d'étude basé sur l'IA: Le PersonaBot. Cet outil permet d'interroger les résultats d'une étude en dialoguant avec des interviewés virtuels, au lieu de consulter des tableaux de résultats.

# IA et données synthétiques

Les données synthétiques sont des données artificielles, c'est-à-dire qui ne proviennent pas directement de l'observation d'individus ou d'événements réels, générées par des modèles statistiques pour reproduire les caractéristiques et les propriétés des données réelles.

On associe souvent les données synthétiques à l'IA mais à l'origine, la synthèse de données est une technique qui permet aux instituts nationaux de statistique de diffuser aux chercheurs des données très fines et riches sur le plan analytique, tout en garantissant la confidentialité et le respect de la vie privée. La théorie autour de la génération de données synthétiques n'est pas récente car elle repose dans ses grands principes sur les méthodes d'imputation. L'imputation de données manquantes, la fusion de données et la modélisation de profils peuvent d'ailleurs s'apparenter à de la génération de données synthétiques.

Les méthodes de génération de données synthétiques ont considérablement évolué avec le développement des méthodes d'apprentissage car elles sont souvent utilisées pour entraîner des modèles de machine learning. L'algorithme Minority **SMOTE** (Synthetic Over-sampling Technique), développé au début des années 2000, permet de réduire les biais des bases d'apprentissage en aioutant des synthétiques pour des groupes insuffisamment représentés. L'IA générative permet aujourd'hui de créer des données synthétiques non structurées.

# Données synthétiques dans les études

L'utilisation des données synthétiques dans les études suscite un certain nombre de fantasmes. Certains y voient une opportunité de réduction du coût des études en remplaçant une partie d'interviewsréellespardesinterviewssynthétiques, d'autres promettent une augmentation de la précision des résultats sur des cibles rares ou difficiles à joindre en augmentant la taille de ces cibles par des interviews synthétiques. Mais ces nouvelles méthodes sont encore loin de faire l'objet d'un consensus théorique dans le domaine des enquêtes. Même si les approches d'imputation ou de fusion sont utilisées depuis longtemps pour compléter ou enrichir des données d'enquêtes ou de panels, le remplacement de répondants réels par des répondants synthétiques soulève encore de nombreuses questions, notamment sur les gains de précision, mais plus globalement sur les protocoles de validation à mettre en œuvre.

Le groupe de travail IA de la commission métier Etudes de Syntec Conseil a publié en juin 2025 une <u>note de position qui clarifie la terminologie autour des données synthétiques et détaille quelques-uns des enjeux méthodologiques, éthiques et juridiques.</u>

Cette note de position s'intègre dans la lignée du « Manifeste pour une utilisation responsable de l'IA dans les sociétés de conseil » de Syntec Conseil dont une déclinaison pour le métier des études a été publiée en mai 2025, endossée par Médiamétrie ainsi qu'une dizaine d'autres sociétés.

# Données synthétiques et mesures hybrides

Selon l'approche retenue, les mesures hybrides peuvent s'apparenter à la création de données synthétiques. C'est le cas notamment des approches de modélisation ou s'appuyant sur des populations synthétiques. En effet, plus le nombre de sources de données à connecter est important, plus la donnée en sortie du processus d'hybridation

sera éloignée des données d'origine. On crée ainsi des individus virtuels qu'on ne peut rapprocher directement d'individus réellement observés.

Le modèle de mesure Cross-Médias de la WFA repose sur ce concept de population synthétique, appelé « Virtual ID model ». La population synthétique est utilisée ici comme socle pour connecter des données individuelles massives provenant de différentes sources et des données issues de panels. L'utilisation de ce modèle de population synthétique permet ainsi à des plateformes et des acteurs médias de partager et de rapprocher des bases de données comportant des données individuelles sensibles, sans compromettre leur confidentialité, et tout en respectant les lois en vigueur.

Les questions soulevées précédemment sur les protocoles de validation à mettre en œuvre restent ouvertes pour ce type d'approches, d'autant qu'on n'a pas toujours de « source of truth » sur un périmètre parfaitement équivalent pour estimer la qualité des résultats de duplication entre les sources obtenues à l'issue du processus. Ceci implique également que, dans un monde où les données sont omniprésentes, les panels et les échantillons demeurent un atout clé. Ils constituent l'unique et indispensable « source of truth » de notre secteur capable de connecter, corriger et dédupliquer les sources de données.

# Enjeux éthiques, transparence et explicabilité

# « L'exigence d'explicabilité devient centrale »

L'IA introduit une nouvelle forme d'opacité méthodologique. Les modèles les plus complexes sont aussi les moins lisibles. Cette complexité pose plusieurs questions:

- Comment auditer un modèle d'IA?
- Comment garantir la non-discrimination dans les estimations ?
- Comment expliquer une méthodologie basée sur de l'IA à des acteurs non techniques ?

L'exigence d'explicabilité devient centrale. Des méthodes dites « white box » (arbres de décision, régressions) sont parfois préférées à des approches plus opaques (deep learning) pour garantir la compréhension du modèle.

Par ailleurs, la **qualité des données d'entrée** reste le socle d'une IA fiable. Un modèle entraîné sur des données biaisées ou déséquilibrées produira des résultats avec les mêmes défauts et pourra même les accentuer.

Enfin, le respect des règles de **protection des données personnelles** (RGPD, anonymisation, consentement) est un prérequis non négociable.

L'intelligence artificielle ne remplace pas la métrologie (méthodes et techniques utilisées pour obtenir une plus grande fiabilité des mesures et assurer leur comparabilité), elle la transforme. Elle ouvre un espace de possibilités inédites pour estimer, anticiper les comportements d'audience.

Mais elle ne dispense pas de rigueur, elle exige même davantage de transparence et de dialogue entre les disciplines : statistique, informatique, sociologie, économie. Son intégration dans les dispositifs de mesure suppose des choix clairs, des arbitrages méthodologiques assumés, et une gouvernance solide.

Dans un monde où les comportements médias se redessinent sans cesse, l'IA est moins une solution magique qu'un nouvel outil au service d'une métrologie adaptative. Un outil puissant, exigeant, et appelé à jouer un rôle clé dans la définition des standards de demain.

« L'IA est moins une solution magique qu'un nouvel outil au service d'une métrologie adaptative »

# IA et accès aux données d'audience

L'IA, et en particulier les LLM, s'apprête à révolutionner l'accès aux résultats d'audience. Actuellement, la consultation de ces données impose l'utilisation d'outils et d'interfaces spécialisés, souvent complexes. L'IA promet de remplacer ces systèmes rigides par une interaction directe et intuitive.

Grâce à des prompts en langage naturel, les professionnels des médias pourront simplement «converser» avec les données pour obtenir des analyses. Cette approche rendra l'exploitation des résultats plus rapide et l'accessible à un plus grand nombre d'utilisateurs.

Le défi majeur reste cependant les modalités de partage de ces données avec ces outils. Aujourd'hui, les acteurs de la mesure d'audience ne partagent pas leurs données propriétaires, ce qui rend les réponses des LLM sur ce sujet peu documentées ou fiables. Le rôle clé de l'IA dans l'univers de la mesure dépendra donc du développement d'un nouveau modèle (technique et économique) permettant aux LLM d'accéder de manière sécurisée et vérifiée à ces bases de données.



# Koenraad Deridder General Manager, CIM

# Réinventer la mesure : Hybride. Modulaire. Cross-média.

Le Centre d'Information sur les Médias (CIM) est un organisme assez unique, puisqu'il est un Joint Industry Committee (JIC) multimédia depuis sa création. Cela nous permet de mutualiser les ressources et de disposer d'une équipe professionnelle dans un petit marché. Cela a également rendu possibles certaines synergies entre les études, par exemple en ayant une seule 'establishment survey'. Ou encore une mesure digitale qui inclut les sites et applications mais aussi, par exemple, les players vidéo, dont les résultats sont depuis 2020 combinés avec les études d'audience TV et d'horodatage afin de fournir un 'rating total vidéo' par programme au quotidien.

Mais une véritable approche cross-média faisait défaut. L'année dernière, nos membres ont décidé que le CIM devait connaître sa « révolution copernicienne » et passer d'une logique centrée sur les médias à une logique centrée sur le consommateur, avec pour objectif de suivre les consommateurs à travers tous les médias et plateformes.

#### Le CIM ONE

Dans ce but, nous avons décidé d'intégrer toutes nos études dans un seul cadre de recherche hybride et modulaire, le **CIM ONE.** Il se compose de cinq modules.

Le **Golden Standard**, déjà publié, décrit la population en termes d'équipements et de consommation média général et sert de référence pour l'univers de l'étude.

Le **ONE panel** rassemble le panel TV et un nouveau panel smartphone, permettant de mesurer autant de consommation média que possible, avec les mêmes personnes, à domicile comme à l'extérieur.

Dans les deux panels, nous mesurons les médias offline via l'audiomatching et tous les médias en ligne, respectivement avec le routeur GfK et le meter Realitymine (par Ipsos).

Le troisième module, le **Data Exchange**, compense les limites statistiques des panels toujours trop petits en intégrant des ensembles de big data. En premier lieu le census internet de Gemius, mais nous avons également mandaté dataBreeders pour mettre en place une solution complète de post-buy vidéo, combinant la télévision linéaire et les données des adserver propriétaires pour la vidéo en ligne afin de produire une évaluation de campagne unifiée.

Toutes ces solutions utilisent un seul ensemble d'outils de modélisation et de data science ainsi qu'une seule plateforme de production. C'est le quatrième module, appelé le **Personification Engine** qui permet d'individualiser et profiler des données de type « appareil ».

Enfin, le cinquième module, est la base de données unique dans laquelle s'effectuent tous les rapports. Nous avons choisi une approche de **population virtuelle**.

#### « A Large Media Model »

Aveccela, nous voulons répondre aux changements profonds dans les attentes des annonceurs, agences et médias au cours des derniers mois. Nos membres s'attendent désormais à ce que nous soutenions leurs approches d'« intelligence prédictive ». Ils doivent pouvoir prévoir ce que les consommateurs verront, entendront ou liront, ce qu'ils en penseront et ressentiront, et comment ils y réagiront. Ces insights doivent être continuellement mis à jour par des « signaux » en temps réel dans leurs « systèmes de connectivité » enrichis par l'IA.

Notre population virtuelle doit donc devenir notre propre 'Large Media Model', en s'intégrant aux 'Large Marketing Model' (comme WPP appelle son système) des membres. Heureusement, avec le CIM ONE, nous disposons des ressources nécessaires.

composants et expertises, grâce à notre orientation multimédia. Lorsque le CIM est né en 1971, de la fusion de deux JIC's de la presse (qui existaient depuis les années 1950), nos fondateurs étaient très visionnaires.

Grâce à notre **Golden Standard**, nous pouvons définirla population de référence pour les modèles. **Les données permettant de prédire ce que les gens voient, entendent ou lisent** ont toujours été notre cœur de métier, et nous resterons la source unique pour les données des médias offline. Nous pensons que nos membres disposeront de leurs propres mesures digitales (outcomes), mais nous pouvons leur fournir des données couvrant l'ensemble du paysage en ligne pour les aider à calibrer et à établir des benchmarks.

Il existe également une demande de la part du marché d'intégrer des **« facteurs qualificatifs »** tels que des scores d'efficacité et de visibilité. Le premier, nous le faisons déjà dans notre planificateur Total Video Advertising (ToVA); le second, dans notre étude out-of-home.

Notre étude OOH nous offre également une expérience en modélisation basée sur l'activité 'Activity Based Modelling' (avec une matrice qui prévoit les 188 millions de déplacements effectués par les Belges au cours d'une semaine moyenne) – un autre pilier des nouveaux systèmes – et une expérience dans la livraison en temps réel de toutes les données via une API CIM dans les systèmes de transaction des membres.

Il semble donc que nous serons en mesure de rester pertinents dans ce nouveau monde enrichi par l'IA, ayant déjà en interne de nombreux

# **Pete Doe**

# Chief Research Officer, Nielsen

# L'utilisation de l'IA par Nielsen pour la mesure de l'audience TV aux États-Unis

L'intelligence artificielle (IA), sous diverses formes, est largement utilisée par Nielsen dans ses produits et processus. Nielsen utilise de grands modèles de langage à travers l'ensemble de l'entreprise, ce qui permet d'améliorer l'efficacité de l'écriture de code, de mieux partager les connaissances et d'optimiser la production de documents internes et destinés aux clients.

L'IA est également intégrée dans nos produits, y compris l'intelligence publicitaire, la mesure de l'audience TV et digitale, la mesure des résultats et la planification. Ce document se concentre sur l'utilisation par Nielsen de l'apprentissage automatique de l'IA dans notre méthodologie de calcul de la mesure Big Data + Panel aux États-Unis. Cette méthodologie intègre 45 millions de foyers « big data » avec le panel de 42 000 foyers/100 000 individus de Nielsen pour fournir ce que nous pensons être la mesure la plus exacte et la plus précise des audiences TV aux États-Unis.

Les «big data» ont une couverture large, mais présentent également des lacunes. Nielsen utilise les données voie de retour (RPD) provenant des décodeurs câble et satellite et les données de reconnaissance automatique de contenu (ACR) provenant des téléviseurs intelligents, et ces deux types de données présentent des défis spécifiques, ainsi que certains défis communs.

Les lacunes sont notamment les suivantes : contenu manquant, TV on/off (en particulier pour les décodeurs), composition et caractéristiques socio-démographiques du foyer, individualisation...

Dans les deux cas, l'identification de foyers et d'appareils communs (foyers «big data» qui sont aussi des foyers de panel) est un élément fondamental de la triangulation du panel et de la data, en utilisant la mesure plus complète du panel comme modèle d'entrée pour combler les lacunes de la data.

Les foyers communs sont identifiés grâce à des mises en correspondance respectueuses de la confidentialité des informations personnelles (noms et adresses, adresse IP) et en utilisant l'apprentissage automatique afin d'identifier les appareils communs dans ces foyers en se basant sur la cohérence des habitudes de visionnage.

Ces foyers et appareils communs, ainsi que d'autres données de panel pertinentes, sont ensuite utilisés pour l'apprentissage du modèle afin de résoudre les problèmes suivants :

## Contenu manquant

Nielsen utilise l'apprentissage automatique sur les foyers et les appareils communs pour reconnaître les schémas associés aux données de visionnage manquantes dans les foyers « big data ». Ce processus quotidien permet d'identifier les foyers susceptibles d'avoir des données de visionnage manquantes et de les exclure des mesures.

#### TV on/off

L'absence d'informations directes sur ce qui est affiché sur les écrans des téléviseurs connectés est un problème bien connu des données RPD provenant de décodeurs. Le téléviseur peut être éteint ou branché sur d'autres entrées alors qu'un boîtier câble reste allumé. Nielsen utilise l'apprentissage automatique sur les foyers communs pour identifier les situations dans lesquelles le décodeur est allumé alors que le téléviseur est sans doute éteint.

## Composition socio-démographique du foyer

Les données RPD et ACR utilisées par Nielsen ne sont pas fournies avec des caractéristiques ou des données socio-démographiques. Ces informations, qui constituent un élément essentiel de la mesure de l'audience TV, sont nécessaires pour déterminer l'audience des individus - le fait de savoir qui vit dans le foyer permet de déterminer qui regarde.

L'utilisation de caractéristiques et de données socio-démographiques attribuées par des tiers est insuffisante : Les études Nielsen ont montré que ces données étaient parfois incomplètes et moins précises. Nielsen utilise ces caractéristiques de tiers, ainsi que les données de visionnage dans les foyers «big data» et les informations de visionnage et les données socio-démographiques des foyers de panels Nielsen, comme entrées dans un réseau de neurones récurrent et une technique

de programmation en nombres entiers mélangés pour identifier les caractéristiques et les données socio-démographiques des foyers «big data».

#### **Individualisation**

L'individualisation des «big

data» constitue une étape essentielle de la mesure. Un élément clé du modèle de Nielsen est l'emplacement de l'appareil dans le logement : un téléviseur installé dans le salon principal aura un public différent de celui d'un téléviseur installé dans une chambre à coucher. Les informations sur l'emplacement du téléviseur sont contenues dans les mesures du panel de Nielsen mais pas dans les données «big data». Nous utilisons donc un modèle d'apprentissage automatique pour attribuer aux données «big data» un emplacement dans le logement, qui est ensuite utilisé comme variable dans le modèle probabiliste d'individualisation qui détermine les

# **Khaled El Serafy**Head of Data Science, BARB

# D'Okner à l'IA: La quête de la mesure hybride

Nous sommes en 1972. Atari vient de sortir Pong, un jeu d'arcade en noir et blanc destiné à être joué sur les téléviseurs monochromes encore répandus en Grande-Bretagne (seule la moitié des ménages britanniques possèdent un téléviseur couleur). IBM vient de lancer la disquette, dont l'usage est pour l'instant strictement limité à l'informatique professionnelle, puisque le premier ordinateur personnel ne verra pas le jour avant plusieurs années.

1972 est également l'année de publication d'un article de Benjamin Okner dans Annals of Economic and Social Measurement : « Constructing a New Data Base from Existing Microdata Sets: The 1966 Merge File. » Dans cet article, Benjamin Okner souligne que l'essor des « ordinateurs électroniques » permettra de répondre à la demande croissante d'informations économiques et démographiques désagrégées.

L'approche d'Okner consiste à utiliser l'« appariement informatique », qui combine deux ensembles de données différents basés sur des variables communes afin de créer un fichier synthétique plus riche. Un premier prototype de ce que nous appelons aujourd'hui la fusion de données.

L'article de Benjamin Okner est l'un des maillons d'une chaîne qui s'étend jusqu'aux systèmes hybrides de mesure de l'audience actuellement déployés ou en cours de construction dans toute l'Europe et au-delà.

#### Brève histoire de la fusion de données

Les méthodologies de fusion de données ont évolué au cours des décennies, nourries par plusieurs contributions issues principalement de publications purement statistiques ou de recherches économiques plutôt que du secteur de la mesure d'audience. Le domaine a progressé depuis les indicateurs de similarité de Mahalanobis (1936), qui sont encore utilisés aujourd'hui dans les fusions de données de mesure d'audience, et le fichier de fusion d'Okner des années 1970, jusqu'aux simulations de Rodgers des années 1980 montrant le risque de biais, et aux méthodes de calibration de Renssen des années 1990. Au cours des dernières décennies, la fusion de données a fait son apparition dans les manuels, les erreurs et les limites ont été clarifiées et l'intégration des big data a fait l'objet d'une attention particulière, dans le respect des contraintes de confidentialité et de représentativité.

Depuis plus de 50 ans, les statisticiens sont aux prises avec les mêmes défis que ceux auxquels nous sommes confrontés aujourd'hui en matière de mesure hybride de l'audience télévisuelle : comment combiner des sources de forces différentes, rapprocher les données au niveau des personnes avec celles au niveau des appareils, et produire des résultats utilisables en cas d'incertitude.

Toutefois, une grande partie de ces recherches ont été effectuées en dehors du monde de la mesure de l'audience télévisuelle, dans les domaines de l'économie, des statistiques officielles, de l'agriculture et de la recherche marketing (un domaine connexe de la mesure d'audience de la télévision). La branche de la fusion de données appliquée à la mesure de l'audience télévisuelle s'est-elle détachée de ses racines originelles ? Nous sommes-nous éloignés des décennies de débats académiques sur la fusion des données ?

Les systèmes hybrides de mesure d'audience n'ont pas tendance à se situer dans cette lignée. Des méthodologies sont proposées, font l'objet d'une validation théorique par des experts du domaine et les résultats sont évalués de manière empirique par des cadres, des chercheurs, des équipes de vente et des statisticiens. En ce sens, bien que la fusion de données soit abordée avec la diligence requise, on observe une distanciation par rapport à la tradition scientifique plus large. Pour que la mesure d'audience fasse de véritables percées à « l'ère de l'IA », nous ne pouvons pas nous contenter d'exiger simplement l'utilisation des dernières technologies lors de la mise au point de nouveaux systèmes hybrides. Il nous faut également réfléchir à la manière dont nous pouvons renouer avec le corpus de recherche sur la fusion de données, qui ne cesse d'évoluer depuis les années 1970.

Deux écoles de pensée en matière de mesure hybride

À ce stade, il convient de noter que la « fusion de données » n'est pas la seule solution en matière de mesure hybride. Il existe deux grandes écoles de pensée lorsqu'il s'agit de combiner les données de panel avec le big data pour produire une meilleure mesure d'audience.

La première école s'apparente à la quête du Graal : un ensemble de données fusionnées unique et polyvalent qui est supérieur à la somme de ses parties à tous points de vue. L'idée est de combiner les données de panel avec des sources provenant des big data (CTV, STB, applications VOD) et de proposer un ensemble de données modélisé qui peut répondre à n'importe quelle question : l'audience des programmes, la couverture multi-plateformes, pourquoi les téléspectateurs changent de chaîne, les mesures d'audience publicitaire de référence, et à peu près tout ce pour quoi nous nous appuyons actuellement sur le panel. L'intérêt est qu'il soit cohérent, que ses résultats soient infiniment flexibles et qu'il combine le meilleur du panel (représentativité, détail au niveau des personnes, étendue) avec le meilleur des big data (profondeur, immédiateté) sans vraiment faire de compromis. Le risque est que ces modèles peuvent être complexes, opaques et difficiles à valider dans tous les cas d'utilisation. Si un élément est erroné, il peut saper la confiance dans l'ensemble du système.

La seconde école de pensée est plus pragmatique. Au lieu d'un ensemble universel de données modélisées, vous concevez des méthodes hybrides distinctes pour des résultats spécifiques. Par exemple, une pour les audiences publicitaires quotidiennes, une autre pour les couvertures des chaînes sur une plus longue période, une autre pour les facteurs d'audience conjointe. Pour tout le reste, on s'appuie sur le panel. L'avantage est qu'il est plus facile de s'y fier : chaque résultat individuel peut être validé de manière indépendante, comparé au panel et audité par des comités sectoriels. L'inconvénient est l'incohérence : l'audience des spots publicitaires d'une source peut ne pas correspondre à l'audience du programme dans lequel ils ont été diffusés selon une autre source. Le champ d'application est également limité, chaque fois que vous avez besoin d'un nouveau rapport, vous devez inventer une nouvelle méthodologie hybride.

Le processus Dovetail Fusion de Barb est un bon exemple de la méthode du Graal. Il combine les données exhaustives de la BVOD avec des données de panel en utilisant une technique de fusion de données pour produire un fichier de données d'audience panel modélisées. CFlight est un bon exemple de la seconde approche plus pragmatique. Les logs des ad-servers sont combinés avec les données du panel, mais ils produisent des rapports très spécifiques sur les performances post-campagne.

#### La fusion de données à l'ère de l'IA

Au fur et à mesure de l'avancée et de l'accessibilité des techniques de traitement et d'analyse des données, la tentation de poursuivre la quête du Graal de la mesure hybride va en s'intensifiant. Imaginons que l'on connecte tous les ensembles de données dont nous disposons : panels, enquêtes, données digitales exhaustives, et que l'on laisse les techniques avancées d'intelligence artificielle produire un ensemble universel de données d'audience, fidèle à ses sources et utilisable pour n'importe quelle analyse. Notre capacité à atteindre cet objectif dépendrait alors de notre tolérance au risque et de nos capacités en tant que data scientists et statisticiens.

Les débats entre les chercheurs sur la fusion des données au cours des cinquante dernières années, en particulier dans le domaine de l'économie, nous incitent à la prudence. L'histoire de la fusion des données relève à la fois de l'évolution et de la révolution. Il ne s'agit pas d'une révolution de la data science et de l'IA dans le sens où il suffirait de connecter les ensembles de données à la dernière technologie et de tourner simplement la manivelle pour obtenir des résultats fiables.

Quelle que soit l'approche que nous adoptons pour aborder la fusion des données, nous sommes confrontés à des problèmes déjà répertoriés dans les travaux de recherche, notamment :

## - Compromis entre biais et variance

Les ensembles de données fusionnés peuvent réduire la variance en s'appuyant sur la force des grands contributeurs, mais potentiellement au risque d'introduire un biais.

# - Erreur d'appariement

Des erreurs peuvent être introduites en fonction de la qualité des « liens de fusion » et des algorithmes d'appariement utilisés.

#### - Quantification de l'incertitude

Nous avons besoin de cadres pour mesurer et décomposer l'erreur introduite par l'appariement, en la distinguant de l'erreur d'échantillonnage ordinaire.

#### - Identification partielle

Dans de nombreux cas, certaines relations entre les variables ne peuvent pas être ré-établies de manière unique ; les méthodes ne peuvent que délimiter ou se rapprocher de l'éventail plausible des résultats.

En se reconnectant aux racines de la recherche sur la fusion des données, il est possible de réduire le risque de réinventer la roue ou de répéter les erreurs du passé.

Voici quelques exemples de ce que l'IA et les techniques de data science apportent désormais. Leur présentation n'a pas pour but de recommander l'utilisation de ces techniques particulières. L'idée est de souligner qu'elles nécessitent une application prudente et adaptée au contexte.

# - Apprentissage de représentations

Les réseaux de neurones et les modèles à variables latentes permettent de découvrir des segments d'audience cachés et des routines comportementales à partir des big data.

#### - Imputation à grande échelle

Les modèles génératifs (VAE, GAN) peuvent être en mesure de compléter des données sociodémographiques manquantes ou d'audience conjointe, en respectant les distributions observées dans les panels.

## - Adaptation de domaine

L'IA peut corriger le décalage de covariables. Par exemple, en transférant des modèles comportementaux construits sur les équipés smart TV vers les foyers équipés TV classique.

#### - Assimilation de données

En s'inspirant des prévisions météorologiques, l'IA peut rapprocher plusieurs « capteurs » bruités (panel, ACR, box) en un modèle d'espace d'état pour la mesure d'audience.

#### - Simulation

Grâce à des populations synthétiques et à des modèles multi-agents, l'IA peut simuler des journaux d'audience quotidiens pour l'ensemble des foyers du pays.

Il est difficile de réussir à sélectionner, combiner et utiliser ces techniques si l'on considère la fusion de données comme une entreprise purement commerciale. Il s'agit également d'un défi scientifique et technique, sur lequel nos collègues de divers secteurs ont travaillé et continuent de travailler depuis plusieurs décennies.

# « Le Graal d'une mesure hybride universelle et polyvalente pourrait bien être atteignable »

Ce livre blanc de Médiamétrie est une étape importante vers la reconnaissance de la mesure hybridecommeunedisciplinescientifique en pleine évolution qui nécessite une collaboration sérieuse au sein du secteur, voire entre différents secteurs. Les initiatives du MRC sur des normes conjointes du secteur pour l'accréditation des solutions de mesure hybride sont également encourageantes. Une collaboration intersectorielle, via le monde de la recherche, avec nos collègues travaillant sur la fusion de données dans le domaine économique, pourrait se révéler d'une grande utilité.

Le Graal d'une mesure hybride universelle et polyvalente pourrait bien être atteignable. En conservant les meilleurs aspects de la mesure par panel et en intégrant les big data en provenance des téléviseurs connectés, des décodeurs et des services de vidéo à la demande/partage de vidéos, elle pourrait être considérée en fin de compte comme une amélioration évidente pour toutes les parties prenantes du secteur. De la même manière que nous considérons que les solutions basées les systèmes audimétriques sont meilleures que celles basées sur les carnets d'écoute dans le contexte du secteur, il se pourrait qu'un jour la mesure hybride soit universellement acceptée comme une version améliorée de la mesure basée uniquement sur le panel. Nous devons la considérer à juste titre comme un défi scientifique et technique majeur que nous pourrons sans doute relever en nous appuyant sur les travaux de recherche déjà établis par nos prédécesseurs.

# **Dierrick Gray**Chief Research & Operations Officer, Numeris

# L'expérience canadienne en matière de mesure hybride (ETAM)

Mesurer les audiences dans le paysage médiatique fragmenté actuel est de plus en plus complexe, mais reste essentiel. Alors que l'audience de la télévision traditionnelle décline, Numeris a introduit au Canada le système Enhanced Television Audience Measurement (ETAM), une solution hybride exploitant les données voie de retour (Return Path Data – RPD) issues des décodeurs (Set-Top Boxes). Cette section résume la mise en œuvre d'ETAM par Numeris, ses avantages, ses défis et son importance pour l'avenir de la mesure d'audience.

# Comprendre la mesure hybride et ETAM

La mesure hybride combine plusieurs méthodes, telles que les données de panel, les données voie de retour (RPD) et les données digitales exhaustives, afin de mieux capturer la consommation média. La solution ETAM de Numeris fusionne les panels traditionnels avec des données RPD à grande échelle pour fournir des insights plus détaillés et précis sur l'audience de la télévision linéaire au Canada.

# Pourquoi la mesure hybride?

La mesure hybride offre des avantages uniques :

- Diffuseurs : Des données plus précises et stables permettent une meilleure compréhension des audiences par chaîne, programme et période.
- Annonceurs : Une meilleure capacité à cibler les campagnes, qu'elles soient larges ou adressables.
- Numeris : Des échantillons plus larges à moindre coût, avec une meilleure intégration, attribution et couverture des chaînes.

# Ce qu'ETAM apporte à l'industrie

- Stabilité améliorée : La mesure hybride ETAM réduit les fluctuations dans les estimations d'audience, produisant des données fiables pour tous niveaux d'audience et favorisant des décisions éclairées.
- Mesure plus précise de l'audience des cibles jeunes : Les méthodes hybrides augmentent les audiences des cibles jeunes l'audience moyenne a augmenté de 8 % (âges 2–17), 4 % (18–34) et 5 % (25–44) montrant que les panels peuvent sous-représenter ces téléspectateurs.
- Réduction des écrans publicitaires à audience nulle : La combinaison des données RPD et des panels traditionnels permet de capter des groupes fragmentés ou de petite taille auparavant manqués en raison des limites liées aux tailles d'échantillon.
- Couverture complète du paysage TV (avec STB): Le panel Numeris repose sur le watermarking audio pour mesurer l'audience TV. Seules les chaînes encodées sont mesurées. Le nombre de chaînes mesurées est passé de 796 à 1 745, augmentant la durée d'écoute totale de 1,4 %. Avec l'ajout des données STB, Numeris suit désormais toutes les chaînes, offrant des insights plus complets sur l'audience.

## Défis liés à la mise en œuvre et à l'utilisation

L'intégration des données panel avec des sources passives à grande échelle comme les données RPD (Return Path Data) est techniquement complexe, nécessitant des modèles avancés et une validation constante pour harmoniser des ensembles de données variés :

- Relation de confiance et gestion des relations: La création de relations de confiance et la gestion des relations avec les nouveaux distributeurs de contenus est devenue essentielle, car Numeris dépend de plusieurs fournisseurs pour produire les résultats d'audience. Ces acteurs n'avaient jamais partagé leurs données auparavant et n'en percevaient pas la valeur. Chaque distributeur nécessite une relation spécifique et une compréhension continue de ses données et technologies.
- Complexité et attentes de l'industrie : L'approche hybride d'ETAM apporte des améliorations mais exige la confiance de l'industrie ; certains s'attendaient à des gains uniformes en volumes de visionnage, mais les résultats ont été mitigés. La complexité de la mesure hybride rend la compréhension et la création de consensus difficiles. L'industrie était habituée à des méthodologies plus simples, ce qui a compliqué la confiance.
- Investissement technologique et expertise nécessaire : La mise en œuvre et la maintenance d'ETAM nécessitent des investissements majeurs en technologie et en compétences. L'évolution constante des données impose flexibilité et adaptation.

« À mesure que les méthodes hybrides deviennent la norme, l'unification des sources de données, la transparence et la communication continue sont cruciales »

Scalabilité du système : Développer un système capable de gérer des ensembles de données volumineux et complexes est crucial mais reste un défi, notamment pour les organisations plus petites.

# L'avenir de la mesure hybride et les enseignements de Numeris

ETAM illustre les avantages et les limites de la mesure hybride de l'audience. L'intégration des big data améliore les insights et s'adapte aux évolutions du paysage médiatique, mais des investissements technologiques précoces et l'implication des clients sont essentiels. Une infrastructure faible a causé des problèmes, et une intégration peu claire a suscité du scepticisme. Un Comité de pilotage composé de clients aurait accéléré l'adoption et amélioré la transparence. À mesure que les méthodes hybrides deviennent la norme, l'unification des sources de données, la transparence et la communication continue sont cruciales. Les solutions hybrides comme ETAM sont indispensables pour des analyses média précises, permettant de relever les défis de fragmentation, complexité et coûts.

# **Dr. Mirko Marr**Head of Research, Mediapulse

# Hi-Res TV en Suisse - Opérationnel et en expansion

Face à un niveau élevé de fragmentation des audiences, le marché suisse de la télévision a demandé à Mediapulse, le JIC TV suisse, de développer une solution permettant d'obtenir des données d'audience plus granulaires et plus stables. En combinant les données de panel avec celles issues des décodeurs (set-top boxes), Mediapulse a livré cette solution, lancée en 2022 et adoptée comme référence pour la planification, le reporting et les transactions. De plus, la plateforme de données hybrides a permis à Mediapulse de mesurer les publicités vidéo digitales sur les CTV et d'intégrer ces données dans la référence TV depuis l'automne 2023.

Le développement de la mesure TV hybride en Suisse, également appelée Hi-Res TV, a été mené par Mediapulse selon une approche « panel-first ». Cela signifie que les données du panel TV fournies par Kantar Media, basées sur un échantillon de 2 000 foyers, sont considérées comme la «source of truth» pour l'audience TV et servent de référence pour entraîner les données issues des décodeurs.

L'entraînement des données des décodeurs inclut l'exclusion des événements invalides, l'estimation des profils des foyers derrière chaque décodeur et l'attribution des audiences aux membres estimés du foyer. Cela aboutit à un panel virtuel basé sur les décodeurs, avec des données d'audience seconde par seconde provenant d'environ 300 000 membres virtuels, dérivées d'environ 150 000 décodeurs (contrats) fournis par les deux principaux distributeurs TV en Suisse.

L'intégration des données du panel virtuel dans les données du panel TV, pour former un panel hybride, repose sur une méthodologie d'imputation et est opérée quotidiennement par Mediapulse. Les données d'audience sont livrées au marché en deux étapes : les données « panel seul » sont publiées le lendemain comme données préliminaires, tandis que les données hybrides deviennent disponibles après trois jours ouvrés et sont considérées comme la référence pour la planification et les transactions. En cas de données manquantes ou invalides provenant des décodeurs, un retour aux données préliminaires pour les jours concernés est prévu.

Comparée aux données « panel seul », la mesure TV hybride en Suisse entraîne une réduction substantielle des écrans publicitaires à audience nulle, une légère augmentation des audiences en journée, une diminution des valeurs aberrantes en termes d'affinité, des données plus stables pour la planification et des analyses d'audience des programmes plus crédibles, notamment pour les chaînes de niche.

La mesure des publicités vidéo pre-roll et mid-roll, diffusées par les ad servers vers les décodeurs, suit la logique de l'approche Hi-Res TV, c'est-à-dire qu'elle combine les données du panel et des décodeurs mais utilise des techniques de mesure différentes. Les contacts publicitaires vidéo sont détectés dans le panel TV via des watermarks cryptées et dans les données des décodeurs via des informations issues des logs.

Les données hybrides pour la TV et les publicités vidéo sont livrées au marché dans un jeu de données intégré, offrant aux annonceurs et aux agences la possibilité de quantifier la couverture et la fréquence dédupliquées pour la planification et le reporting des campagnes TV/vidéo crossmedia.

# Andrea Mezzasalma Founder & CEO, dataBreeders

# Mesure hybride des médias à l'ère de l'IA : Opportunités et limites

Tout d'abord, une clarification sémantique du terme d'« IA ». Dans une interprétation large, l'IA peut inclure de nombreuses méthodes avancées d'apprentissage automatique, y compris des techniques non neuronales telles que la méthode des k plus proches voisins (KNN), les forêts aléatoires et le gradient boosting. Toutefois, aux fins du présent document, nous limitons l'«IA » à l'apprentissage profond (réseaux neuronaux avec de nombreux nœuds et couches), aux grands modèles de langage (LLM), à la conversion de la parole en texte, à la vision par ordinateur, à l'IA générative et aux approches connexes.

Chez dataBreeders, nous utilisons déjà largement les LLM pour améliorer la productivité - par exemple pour le codage, la découverte, l'assistance aux utilisateurs, la documentation, la traduction, l'amélioration des textes, etc.

Dans le domaine de la mesure hybride, notre position actuelle n'est pas de remplacer nos modèles de base par l'IA. Nous visons plutôt à exploiter l'IA avancée pour améliorer l'infrastructure et le flux de travail de bout en bout, tout en préservant la robustesse de notre modélisation de base.

#### Modèles de mesure hybrides de base

Pour nos modèles hybrides de base, nous nous appuyons actuellement sur un mélange de méthodes non neuronales, notamment :

- Modélisation statistique basée sur des distributions bien connues (parexemple, processus de Poisson composé, binomial composé), entraînée par des méthodes d'apprentissage automatique telles que la validation croisée k-fold.

- Modèles linéaires généralisés (par exemple loglinéaire, logit, etc.)
- KNN (k plus proches voisins)
- Gradient boosting

Au fil du temps, ces modèles ont démontré plusieurs avantages par rapport aux méthodes d'IA « avancées » :

- De bonnes performances, même sur des ensembles de données d'apprentissage relativement petits (par exemple, des « panels de référence à source unique »)
- Transparence et auditabilité accrues, contrairement aux modèles « boîte noire » tels que les LLM pré-entraînés
- Capacité à exploiter les relations causales sousjacentes qui sont plus susceptibles de persister dans le temps, plutôt que des modèles purement associatifs qui peuvent être fallacieux ou changer au cours des périodes ultérieures.
- Prévisibilité et contrôlabilité : nous pouvons estimer l'exactitude et la précision (« biais » et « variance » dans le terminologie de l'apprentissage automatique), ainsi que la régression à la moyenne, et mettre en œuvre avec précision des choix spécifiques pour améliorer l'une ou l'autre de ces qualités, ou trouver le meilleur équilibre entre elles.

Compte tenu de tout cela, nous n'avons pas l'intention, à court terme, de remplacer ces outils par des solutions d'apprentissage profond (qui, après tout, existent depuis des décennies) ou des LLM pré-entraînés.

Nous restons toutefois ouverts à une réévaluation. Si de nouvelles exigences apparaissent et que nos outils actuels ne peuvent les satisfaire, ou si les méthodes d'IA deviennent suffisamment fiables et transparentes, nous pourrons reconsidérer la question.

#### Détection et amélioration des attributs

Nous nous concentrons plutôt sur l'utilisation de l'IA pour la détection et l'amélioration des caractéristiques. Aucune sortie de modèle ne peut dépasser la qualité de ses données d'entrée et, dans des contextes tels que la mesure des médias, l'IA peut grandement aider à résumer et à classer le contenu (texte, synopsis, audio, vidéo) en genres, ambiances, orientations, etc.

Par exemple, si nous pouvons classer avec plus de précision le contenu visionné sur le téléviseur d'une famille, nous pouvons mieux estimer la probabilité que les différents membres de cette famille le regardent et en déduire leurs profils démographiques.

Dans ce domaine, l'IA surpasse largement l'homme en termes d'évolutivité et de cohérence entre les plateformes. En conséquence, nous expérimentons activement des modules pilotés par l'IA pour compléter notre pipeline de mesures et améliorer la qualité globale.

#### Répondants synthétiques

Une autre application prometteuse de l'IA générative est la construction de répondants synthétiques, c'est-à-dire la génération de microdonnées qui ressemblent à des ensembles de données réels tout en minimisant le risque de divulgation d'informations sur des individus réels.

Toutefois, de manière générale, nous restons prudents à l'égard des répondants synthétiques. Dans la mesure du possible, nous préférons conserver les données à caractère personnel dans des environnements protégés et séparés et fournir des rapports à la demande via des API ou des outils de reporting. En fait, nous partons du principe que la génération de répondants synthétiques dégradera inévitablement certaines des corrélations complexes présentes dans les

données originales et, de toute évidence, les répondants synthétiques ne peuvent pas être utilisés pour le ciblage et l'optimisation.

Néanmoins, pour être prêts à répondre aux exigences futures des clients, nous expérimentons de manière proactive des modèles tels que les GAN, les VAE et les LLM, ainsi que des stratégies plus traditionnelles telles que l'omission aléatoire d'une partie des données d'origine et l'inférence de données de remplacement.

Il convient toutefois de noter que nous ne considérons pas les répondants synthétiques comme une méthode de modélisation centrale dans la mesure hybride. Nous les considérons plutôt comme une option tactique (appliquée en amont ou en aval du pipeline de modélisation de base) pour garantir la conformité avec les exigences en matière de protection de la confidentialité, si nécessaire.

# Donner aux utilisateurs les moyens d'agir grâce à des requêtes en langage naturel

Enfin, nous développons des outils de reporting basés sur des LLM qui utilisent des invites en langage naturel pour interagir avec nos API et nos points de terminaison, préservant ainsi le principe du calcul multipartite dans des environnements séparés afin de garantir une sécurité maximale des données.

# **Zuber Nosimohomed**

# Chief Product and Business Officer - Kantar Media / TechEdge - Président

# Deux méthodes différentes mais complémentaires : la combinaison des panels et des données numériques permet d'obtenir de meilleurs résultats

La consommation du public est de plus en plus fragmentée entre les appareils et les plateformes. Cela ne devrait pas être le cas pour la mesure. La mesure hybride - qui combine les panels et les données numériques - génère des données claires, comparables, précises et transparentes. En d'autres termes, ils sont complémentaires :

- Les panels observent directement ce qu'un échantillon national représentatif de personnes réelles regarde - sur n'importe quel appareil d'une manière qui ne dépend pas de la source du média et qui ne porte pas atteinte à la vie privée.
- Les données numériques sont des informations granulaires à grande échelle qui permettent de savoir ce qui est regardé sur chaque appareil.

Sans données numériques, il est difficile de mesurer avec précision tout ce que les audiences regardent. D'autre part, les données numériques qui ne reflètent pas une « vérité de terrain » basée sur les personnes seront toujours sujettes à des biais, manqueront de comparabilité et d'indépendance... la modélisation a ses limites.

# Fonctionnement - les méthodes de fusion et de calibration sont les ingrédients clés de la « recette »

Un bon système de mesure hybride doit combiner les meilleures méthodologies, des technologies respectueuses de la confidentialité et des outils clients simples et faciles à utiliser. Bien qu'un certain nombre de ces éléments puissent avoir une dimension mondiale, il est essentiel de toujours refléter les besoins du marché local. Ces systèmes de mesure adopteront l'une des

deux approches suivantes (ou une combinaison des deux) :

- Fusion fusion d'ensembles de données pour créer des données détaillées, y compris des données numériques et des données mesurées par des panels.
- Calibration ajustement d'ensembles de données numériques déjà détaillées afin de refléter une « vérité de terrain » indépendante basée sur les personnes, mesurée par le panel.

Dans les deux cas, ces approches peuvent être améliorées par des données synthétiques et des technologies de protection de la confidentialité afin de garantir leur précision tout en protégeant la confidentialité de l'utilisateur.

En outre, les solutions hybrides n'ont de valeur que si les données sont facilement accessibles et utilisables par les clients. Les outils d'analyse doivent donc être en mesure de traiter de grands volumes de données pour fournir aux utilisateurs des informations utiles en temps voulu. Ce n'est pas encore le cas pour de nombreux outils disponibles sur le marché aujourd'hui, ce qui peut constituer un obstacle à l'adoption rapide de la mesure hybride.

Localisation - exemples de la manière dont les approches de mesure hybride sont adaptées avec succès aux besoins spécifiques des marchés locaux.

**Pays-Bas :** NMO a une vision audacieuse de la mesure cross-média et vise à combiner la télévision, la vidéo en ligne, l'audio, la presse et, à l'avenir, les données OOH. Grâce à ce partenariat, nous combinons nos panels TV et nos données

BVOD, et travaillons en étroite collaboration avec NMO, Ipsos et d'autres pour fournir une vue d'ensemble des audiences avec une déduplication précise.

**Norvège :** En Norvège, nous travaillons avec le Comité des propriétaires de médias (Media Owners Committee) pour fusionner deux panels complémentaires - à domicile (in-home) et hors domicile (out-of-home) - et calibrer les données de recensement BVOD afin de fournir une image unifiée et quotidienne des audiences de la radiodiffusion et de la diffusion en continu, tant pour le contenu que pour les publicités, à domicile et hors domicile.

**Brésil :** Notre outil Cross-Media Planner intègre et calibre les données de panel avec les données de YouTube, fournissant ainsi aux agences et aux annonceurs un outil de planification de nouvelle génération, fondé sur un panel de source unique et pas seulement sur une modélisation probabiliste, avec la granularité de données dynamiques en temps réel.

Royaume-Uni: Au Royaume-Uni, nous travaillons en partenariat avec Barb pour calibrer les données de panel et de recensement BVOD en utilisant l'expansion virtuelle pour donner une vue unifiée de la consommation de contenu radiodiffusé sur tous les écrans. La prochaine initiative de Barb avec Panel Plus consiste à fusionner les données des décodeurs et des téléviseurs connectés afin de fournir des mesures de l'audience TV plus précises.

**Royaume-Uni et États-Unis :** Pour le compte des programmes Origin et Aquila, nous intégrons les données des panels aux expositions

« L'utilisation de l'IA de manière transparente, explicable et vérifiable (...) rendra les systèmes de mesure hybride plus efficaces et plus précis » publicitaires des plateformes numériques (d'autres ensembles de données et types de médias sont prévus) afin de fournir des résultats post-campagne par le biais d'une méthode unique, respectueuse de la confidentialité, qui s'appuie sur la « modélisation de l'identification virtuelle ».

# IA - Au lieu de les remplacer, l'IA complète et améliore les solutions de mesure hybride

L'IA et l'analyse avancée transforment les systèmes de mesure de bout en bout, y compris (i) l'ensemble du cycle de vie des panels; (ii) les développements de la data science en imputant des données pour combler les lacunes, en développant des outils prédictifs et d'autres approches de pointe; (iii) dans les outils des clients finaux pour simplifier massivement la façon d'extraire des informations à partir d'ensembles de données complexes.

L'utilisation de l'IA de manière transparente, explicable et vérifiable, dans le respect des nouvelles réglementations émergentes, rendra les systèmes de mesure hybride plus efficaces et plus précis. Elle peut également faire passer la protection de la confidentialité à un niveau supérieur, en permettant l'utilisation de nouvelles techniques qui évitent d'exposer des données individuelles et résolvent des problèmes cruciaux entre les canaux, comme les différences de taxonomie. Toutefois, l'IA repose sur des données d'entrée de haute qualité provenant de sources numériques et de panels. Elle aura donc davantage tendance à compléter, plutôt qu'à remplacer, les solutions de données hybrides.

#### En résumé...

La mesure hybride, qui combine les panels et les données numériques, offre le meilleur des deux mondes en fournissant l'indépendance, la transparence, la précision et la comparabilité dont les diffuseurs, les annonceurs, les agences et les plateformes ont besoin. En combinant les bonnes données, dans le bon contexte et de la bonne manière, la mesure hybride peut apporter de la clarté à un paysage médiatique complexe. Or, c'est précisément ce dont a besoin le secteur de la mesure...

# **Mario Paic**

# Chief Research Officer, Audience Measurement, Ipsos

Comme le souligne ce livre blanc, méthodologies de base du secteur de mesure d'audience ne cessent d'innover d'évoluer. Ces approches avancées ont depuis longtemps dépassé le stade de l'observation d'échantillons de comportements individuels par le biais d'enquêtes ou de panels et commencent à intégrer une multitude de sources de données. L'intégration des données a toujours été composante cruciale de la mesure de l'audience : qu'il s'agisse des premières tentatives d'intégration de données de consommation produits dans les mesures TV de référence afin d'améliorer les capacités de ciblage et d'efficacité ou dans le cadre des premières fusions de données de référence conçues pour mesurer la couverture cross-plateformes des marques de médias.

Aujourd'hui, un nombre croissant de standards de mesure d'audience utilisent des cadres méthodologiques hybrides entièrement développés, en particulier ceux qui fusionnent des données de panel avec des ensembles de données complets de recensements, collectés par le biais de technologies RPD ou ACR, de tags third-party et de SDK, etc.

Aufildutemps, ces méthodes ont considérablement bénéficié des progrès de la data science ainsi que de la disponibilité accrue du cloud computing, ce qui a permis de mettre en place des approches plus sophistiquées, tirant parti des volumes croissants de données collectées par le biais d'un large éventail de technologies.

Ipsos a été à l'avant-garde de ces développements et a mis au point certaines des premières méthodologies hybrides approuvées par le secteur avec son travail de pionnier sur les standards de mesure d'audience hors domicile (OOH) dans le monde entier. La création d'un standard OOH implique la tâche complexe de combiner un grand nombre de sources de données (par exemple, des enquêtes sur les déplacements s'appuyant sur les données GPS, des comptages de trafic, des données cartographiques numérisées, etc.) Cette complexité pose des défis uniques, notamment en ce qui concerne l'intégration et la modélisation de toutes les données en un ensemble cohérent, tout en préservant l'intégrité statistique de chaque source de données.

Fort de ces avancées, nous avons encore affiné nos méthodologies en nous concentrant récemment sur l'exploitation de données synthétiques.

Dans le contexte actuel, le terme de « données synthétiques » peut recouvrir plusieurs définitions. Il est important de souligner que l'utilisation de données synthétiques dans la mesure de l'audience diffère nettement des cas d'utilisation récents dans le secteur plus large des études de marché. Plus précisément, elle diffère des applications qui utilisent l'IA générative pour extrapoler des informations à partir de données synthétiques.

Nous utilisons des données synthétiques sous la forme d'une population synthétique représentative respectueuse de la confidentialité, qui sert de couche d'intégration reliant les données de panel, les signaux des recensements et les intégrations serveur à serveur. Différentes techniques d'intégration des données peuvent être employées pour intégrer les ensembles de données, en fonction du cas d'utilisation et du type de données (par exemple, données agrégées ou désagrégées). Des couches de pondération et de calibration garantissent la fidélité, tandis que des contrôles de qualité rigoureux et la

« La compréhension du comportement des individus par l'observation directe en utilisant des méthodes éprouvées telles que les panels représentatifs de haute qualité et des enquêtes, continuera à jouer un rôle essentiel dans la construction de systèmes de mesure d'audience fiables à l'avenir »

transparence permettent d'attester la fiabilité de la méthodologie. Il est essentiel que ces ensembles de données soient conçus de manière à empêcher toute réidentification, afin de respecter les obligations en matière de protection de la confidentialité sans pour autant sacrifier l'utilité analytique. Tout cela nous permet de produire des couvertures et répétitions dédupliquées et granulaires, des rapports cohérents via des API et des interfaces utilisateur, et des standards crédibles que les clients peuvent utiliser pour planifier, évaluer et activer.

Le meilleur exemple d'utilisation de cette approche est Ipsos iris, la mesure digitale de référence au Royaume-Uni et en Australie, où les données d'un panel automatique single-source sont combinées avec diverses sources de données au sein d'une population synthétique représentative.

L'adoption par le secteur, et l'acceptation par le marché, de données de référence qui combinent les mesures indépendantes traditionnelles avec diverses autres sources de données, y compris celles collectées et fournies directement par les propriétaires de médias et les plateformes, produites grâce à des techniques avancées de modélisation et d'intégration de données, y compris l'utilisation de données synthétiques, démontrent que la transformation des systèmes de mesure de l'audience est déjà en cours. Il est probable que cette transformation devra se poursuivre, voire s'accélérer, en raison des progrès de l'IA et de son impact sur la manière dont le contenu et la publicité sont produits, distribués et consommés.

Cela dit, chez Ipsos, nous pensons que la compréhension du comportement des individus par l'observation directe en utilisant des méthodes éprouvées telles que les panels représentatifs de haute qualité et des enquêtes, continuera à jouer un rôle essentiel dans la construction de systèmes de mesure d'audience fiables à l'avenir, non pas malgré mais en raison des progrès de l'IA et de son impact potentiel sur le secteur.

# Synthèse et prochaines étapes: Hybride & lA au service des nouveaux enjeux de la mesure

# Synthèse et next steps : **Hybride & IA** au service des nouveaux enjeux de la mesure

Dans un contexte de délinéarisation et de fragmentation croissantes des usages, les frontières entre diffuseurs linéaires et plateformes de streaming s'estompent et les logiques de convergence s'accélèrent.

Cette transformation pour les mesureurs d'audience se traduit par une complexité croissante et une nécessité d'adaptabilité des protocoles de mesure.

# « L'hybridation n'est pas un concept théorique »

Face à la digitalisation des médias, l'arrivée de nouveaux acteurs et la diversification des modèles de distribution, la mesure d'audience doit intégrer une composante data croissante pour répondre aux défis liés à la précision insuffisante des panels pour estimer la consommation d'usages toujours plus fragmentés.

L'hybridation des données panel avec des données voie de retour - des data - nous semble donc une évolution inéluctable et généralisée. L'hybridation n'est pas un concept théorique. L'avenir de la mesure d'audience est déjà là et il repose sur ces nouvelles méthodologies hybrides. La très grande majorité des nouvelles mesures publicitaires et éditoriales du Total Video s'appuient à la fois sur des données de panel/échantillon et sur des Big Data. Le framework HALO de la WFA ou la nouvelle mesure Big Data TV de Nielsen sont de claires illustrations de cette tendance.

Comme nous l'avons exposé dans ce document, l'hybride doit permettre à la mesure de relever trois défis clés : étendre la couverture de la mesure, répondre à la fragmentation des médias avec des résultats plus granulaires à un coût raisonnable, et fournir des résultats dans un écosystème totalement nouveau (les plateformes d'AdTech en particulier).

La construction de ces mesures nécessite plusieurs ensembles de données pour fonctionner. L'accès à ces données soulève de multiples questions. Un nombre croissant de réglementations des médias et du numérique abordent cette question et poussent en faveur d'une transparence / partage des données.

Les comités interprofessionnels (JIC) et les sociétés de mesure d'audience sont bien équipés pour créer les normes nécessaires et auditer ces données. Mais une nouvelle ère de collaboration doit émerger. Les systèmes de mesure d'audience ont été construits comme des silos indépendants, d'un média à l'autre et d'un pays à l'autre. Si la mesure des médias doit rester locale et conserver une forme de souveraineté, il existe un besoin évident de les connecter à une vue plus globale. Les flux de données provenant des plateformes sont mondiaux. Les annonceurs souhaitent aussi obtenir une vue globale. Et avec la croissance du numérique, la réglementation est aussi de plus en plus mondiale (en particulier en Europe). Cela implique certainement qu'une plus grande collaboration s'instaure entre les acteurs de la mesure. Un mouvement qui est clairement en cours. La coalition pour la mesure d'audience (AMC), qui rassemble de nombreux acteurs de notre industrie, en est un signal clair.

# « Les panels sont centraux car ils ont pour rôle de faire le pont entre les différentes données »

Nous avons tendance à décrire le système hybride comme un système dual combinant des ensembles de données et des panels. Si c'est factuellement vrai, le panel doit être considéré comme la pièce centrale. Pas nécessairement plus important, mais central. Sans panel, nous ne serions pas en mesure de construire une mesure. Sans ensembles de données, nous construirions une mesure avec une résolution plus faible. Les panels sont centraux car ils ont pour rôle de faire le pont entre les différentes données, soit directement, soit à travers les modèles qu'ils permettront de construire. Les panels fournissent les données neutres et indépendantes nécessaires pour mesurer les données individuelles et la duplication. Comme l'affirme la WFA, ils sont une « source de vérité » indépendante et neutre. La neutralité et l'indépendance sont bien sûr essentielles pour instaurer la confiance. La duplication est nécessaire, car c'est la valeur clé qui permet de rassembler les divers ensembles de données et de compter le nombre d'individus qui ont vu un contenu ou une publicité. Les panels sont bien sûr confrontés à des défis, car ils deviennent plus difficiles à recruter et plus coûteux. Mais investir en eux est une nécessité absolue pour fournir la prochaine génération de mesure hybride.

Depuis le lancement de ChatGPT 3.5, l'IA émerge à grand pas et prépare une nouvelle révolution numérique, accélérant les tendances de fragmentation des usages déjà observées. L'IA promet aussi de nouvelles perspectives de transformation des mesures. Mais de même que les data n'ont pas remplacé les panels, l'IA ne sera pas la révolution que certains attendent. Les panels synthétiques ne sont pas un substitut aux échantillons traditionnels. Ce sont de nouveaux modèles qui nous aident à faire progresser des méthodologies hybrides, en utilisant à la fois les échantillons comme source de vérité et des data. Le futur sera définitivement hybride.

« Mais de même que les data n'ont pas remplacé les panels, l'IA ne sera pas la révolution que certains attendent (...) Le futur sera définitivement hybride »

# Liens et bibliographies

#### **MRNews:**

(2025) Podcast #7 - Données synthétiques : quelle place dans les études marketing ? Avec Stéphane Marcel, Président de Syntec Etudes

## Notes de position de la commission Etudes de Syntec Conseil :

(2025) Manifeste pour une utilisation responsable de l'IA dans les sociétés de conseil (2025) Données synthétiques et études marketing & sondages d'opinion

#### **ESOMAR Practical Guidance:**

(2024) 20 questions to help buyers of AI-based services for market research and insights (2025) 5 Topics of Discussion to Help Buyers of Augmented Synthetic Data

## **UKOM White Paper:**

(2024) Synthetic Data Gets Real

# **Ipsos Papers on AI and market research:**

(2023) Conversations with AI Part I: How generative AI and qualitative research will benefit each other

(2023) Conversations with AI Part II: Unveiling AI quality in qualitative workstreams

(2023) Conversations with AI Part III: How AI boosts human creativity in ideation workshops

(2024) Conversations with AI Part IV: AI-assisted knowledge libraries and curation, the search for a trusted output

(2024) Conversations with AI Part V: Is there depth and empathy with AI twins?

(2024) Conversations with AI Part VI: AI-powered moderator bots: Enhancing empathy or eroding connection?

#### Publications de l'EGTA sur la mesure d'audience :

Radio & Audio

TV

## Comscore:

(2009) An early take on the topic of Hybrid audience measurement

# Articles et conference papers publiés par les équipes scientifiques de Médiamétrie :

(2022) <u>Fusion de fichiers d'enquêtes</u>, Vanheuverzwyn A., dans Données manquantes, Éditions Technip, Paris, p. 165-187.

**(2019) Learning on Survey Data to Qualify Big Data: application to web data**, Vanheuverzwyn A., 11th MAASC Day, Besançon.

**(2019) Mixed Use of Big Data and Survey Data for Media Audience Measurement in France : An** <u>Overview</u>, Dudoignon L., Le Sager F. et Vanheuverzwyn A., 8th European Survey Research Association Conference, Zagreb.

**(2018) Learning on Survey Data to Qualify Big Data in a Web Environment**, Duprat L., BigSurv 2018, Barcelona.

(2018) Big data et mesure d'audience : un mariage de raison ?, Dudoignon L., Le Sager F. et Vanheuverzwyn A., Economie et Statistique / Economics and Statistics, 505-506, p. 133-146.

(2016) Mesure d'audience et données massives : mythes et réalités, Vanheuverzwyn A., 9ème Colloque

Francophone sur les Sondages, Gatineau.

(2014) Mesure hybride de l'audience TV, Dudoignon, L. et Logeart J.. 46èmes Journées de Statistique de la SFdS, Rennes.

(2013) Panorama des dispositifs hybrides de mesure d'audience des médias en France, Vanheuverzwyn A. et Zydorczak L., Symposium international de Statistique Canada sur les questions de méthodologie, Ottawa.

**(2012) Enquêtes et données exhaustives : un nouveau défi pour les mesures d'audience**, Dudoignon L. et Zydorczak L., 7ème Colloque Francophone sur les Sondages, Rennes.

**(2012) Médiaplanning & hybridation appliqués à l'Internet mobile**, Crochet G. et Santini G., 44èmes Journées de Statistique de la SFdS, Bruxelles.

(2011) Méthode d'hybridation de données appliquée à la mesure d'audience de l'Internet Mobile en France, Vanheuverzwyn A. et Vouge E., dans les Actes des 43èmes Journées de Statistique de la SFdS, 27 mai 2011, p. 158-163.

# Les auteurs



**Aurélie Vanheuverzwyn,**Directrice Exécutive - Data et Méthodes

Diplômée de l'École nationale de la statistique et de l'analyse de l'information (ENSAI), Aurélie Vanheuverzwyn a rejoint Médiamétrie en 1999. Actuellement Directrice Exécutive - Data et Méthodes de Médiamétrie, elle pilote les activités Data Science et coordonne les fonctions de Data Privacy et Information Security Officer. Aurélie Vanheuverzwyn est également membre élue de l'International Statistical Institute (ISI) et de la Société Française de Statistique (SFdS).



**Julien Rosanvallon,** *Directeur Général Adjoint - Marketing et Expérience Client* 

Diplômé de la Maîtrise des Sciences et Gestion de Dauphine et d'un Master à l'IEP de Paris, Julien Rosanvallon débute sa carrière au sein de Orange Advertising comme Responsable des Etudes. Il intègre Médiamétrie en 2003 en tant que Directeur de la filiale Médiamétrie//NetRatings. Depuis 2020, il est Directeur Général Adjoint en charge des Mesures d'Audience. Il est l'auteur de trois ouvrages sur les médias, la mesure et le numérique dont «La Mesure des Médias» paru en 2025 aux Editions Economica.