

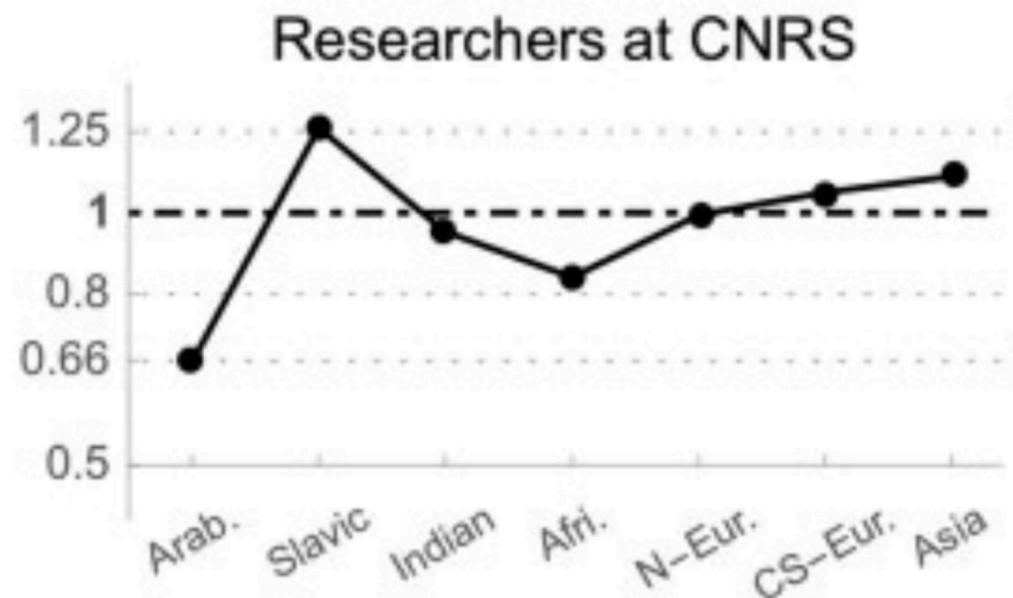
Appréhender les discriminations à l'aide de l'analyse d'image et de l'onomastique

Antoine Mazieres

Le 6 novembre 2020, Séminaire *Systemes Complexes en Sciences Sociales*, LABEX SMS.

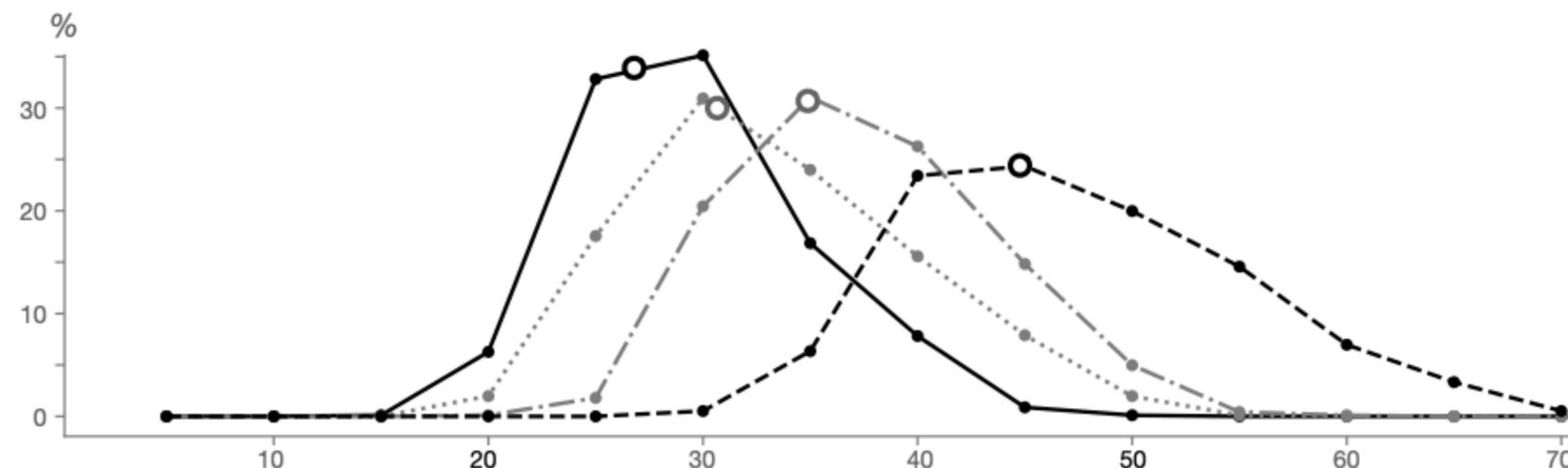
Présentation de deux articles

- Mazieres, Antoine et Roth, Camille (2018). **Large-scale diversity estimation through surname origine inference.** *Bulletin of Sociological Methodology.*



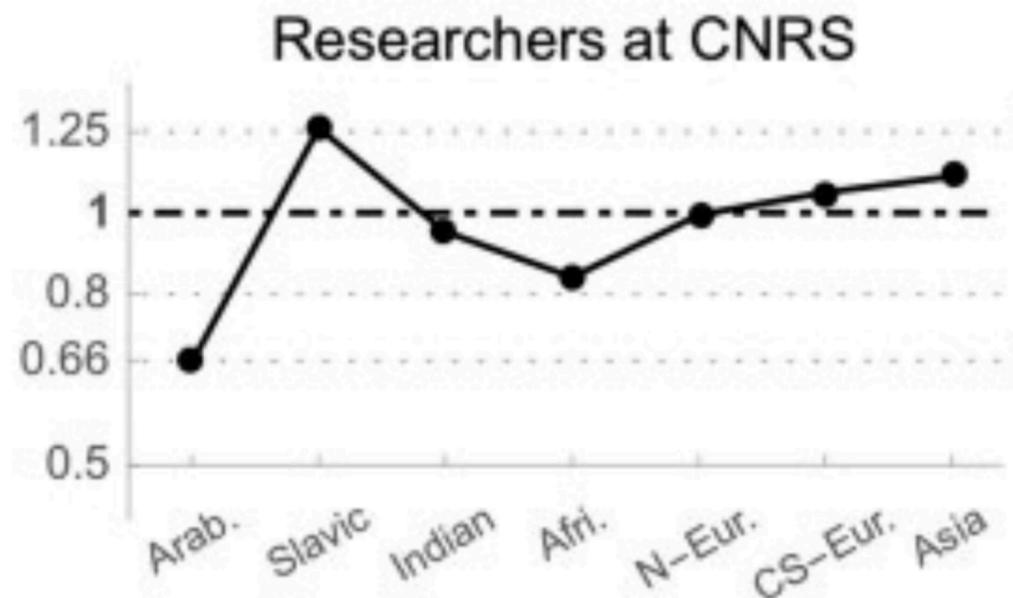
- Mazieres, Antoine, Menezes, Telmo et Roth, Camille (2020). **Computational appraisal of gender representativeness in popular movies.** *Arxiv.*

Evolution of female face ratio in popular movies



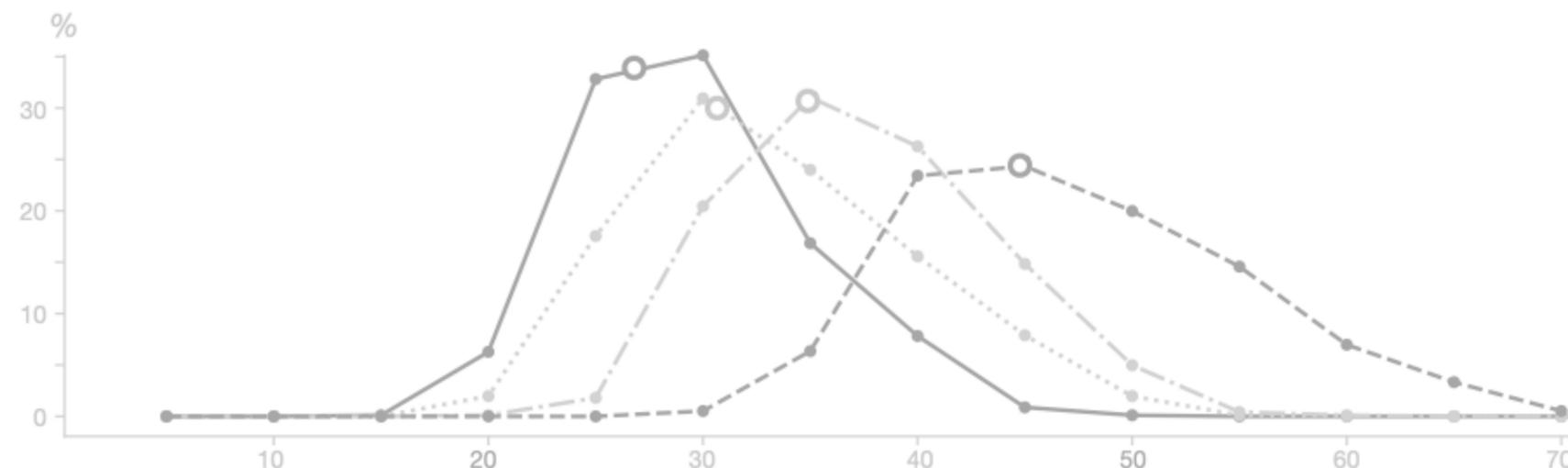
Présentation de deux articles

- Mazieres, Antoine et Roth, Camille (2018). **Large-scale diversity estimation through surname origine inference.** *Bulletin of Sociological Methodology.*



- Mazieres, Antoine, Menezes, Telmo et Roth, Camille (2020). **Computational appraisal of gender representativeness in popular movies.** *Arxiv.*

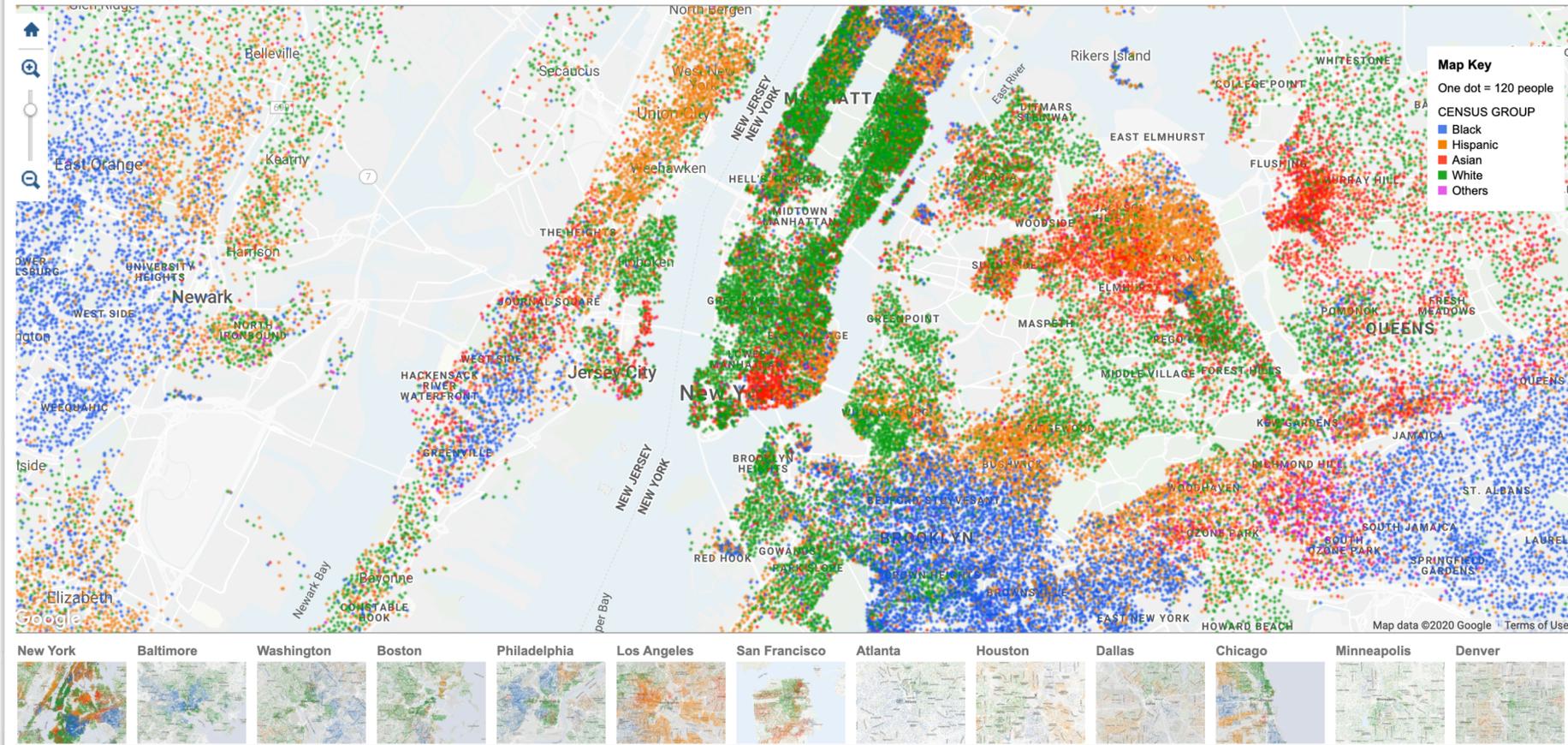
Evolution of female face ratio in popular movies



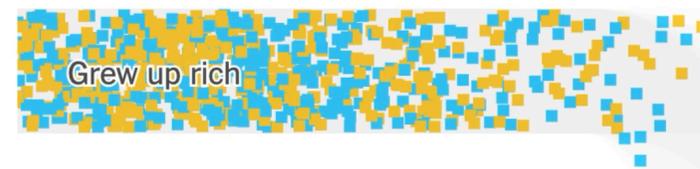
Pourquoi mesurer les discriminations par origine en France ?

Mapping Segregation

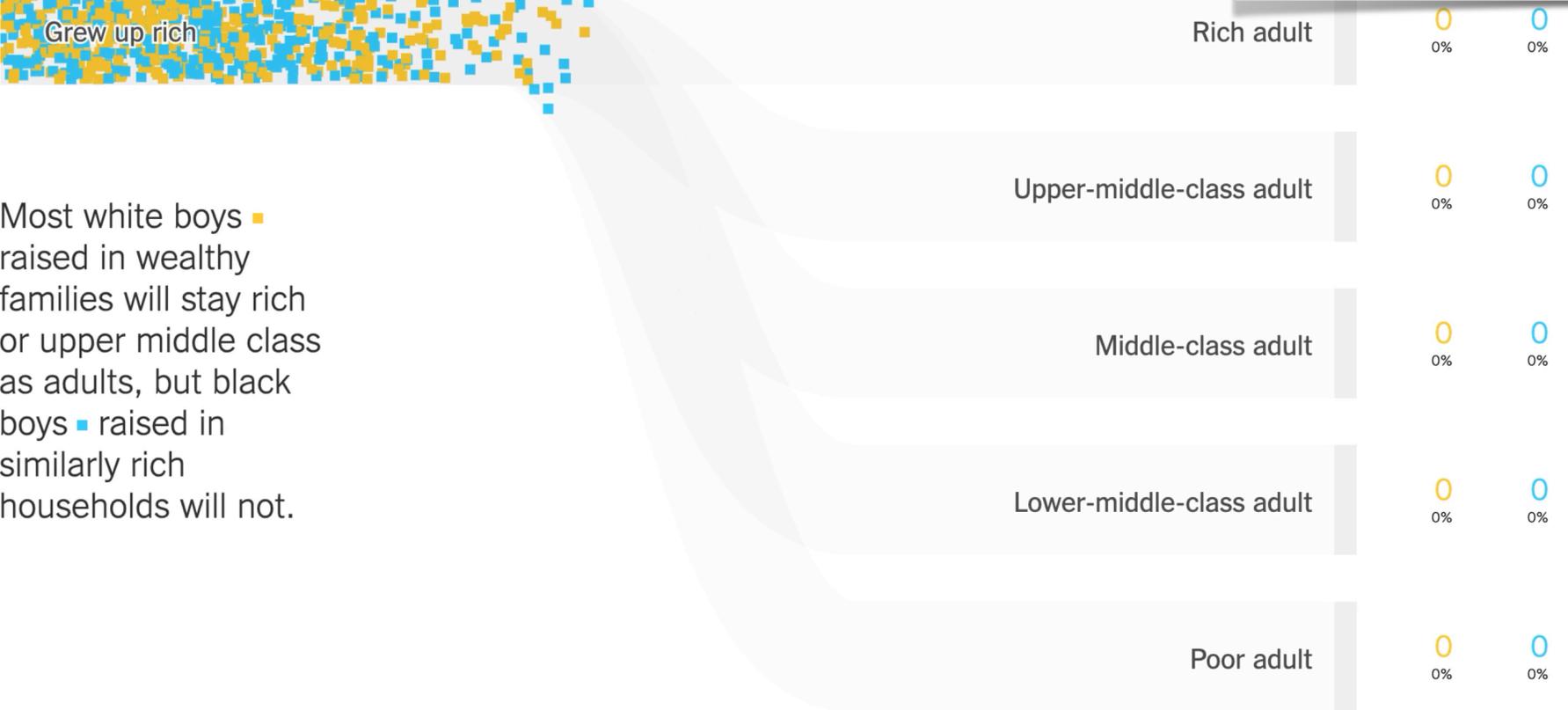
New government rules will require all cities and towns receiving federal housing funds to assess patterns of segregation.



Follow the lives of 792 boys who grew up in rich families ...



Most white boys raised in wealthy families will stay rich or upper middle class as adults, but black boys raised in similarly rich households will not.



Adult outcomes reflect household incomes in 2014 and 2015.

Racismo no Brasil

Impacto



Indicadores	Brasileiro branco	Brasileiro negro
Analfabetismo ^[77]	5,9%	13,3%
Nível universitário ^[78]	15,0%	4,7%
Expectativa de vida ^[79]	73,13	67,03
Desemprego ^[80]	5,7%	7,1%
PIB per capita ^[81]	R\$ 22,699	R\$ 15,068
Mortes por homicídios ^[82]	29,24%	64,09%

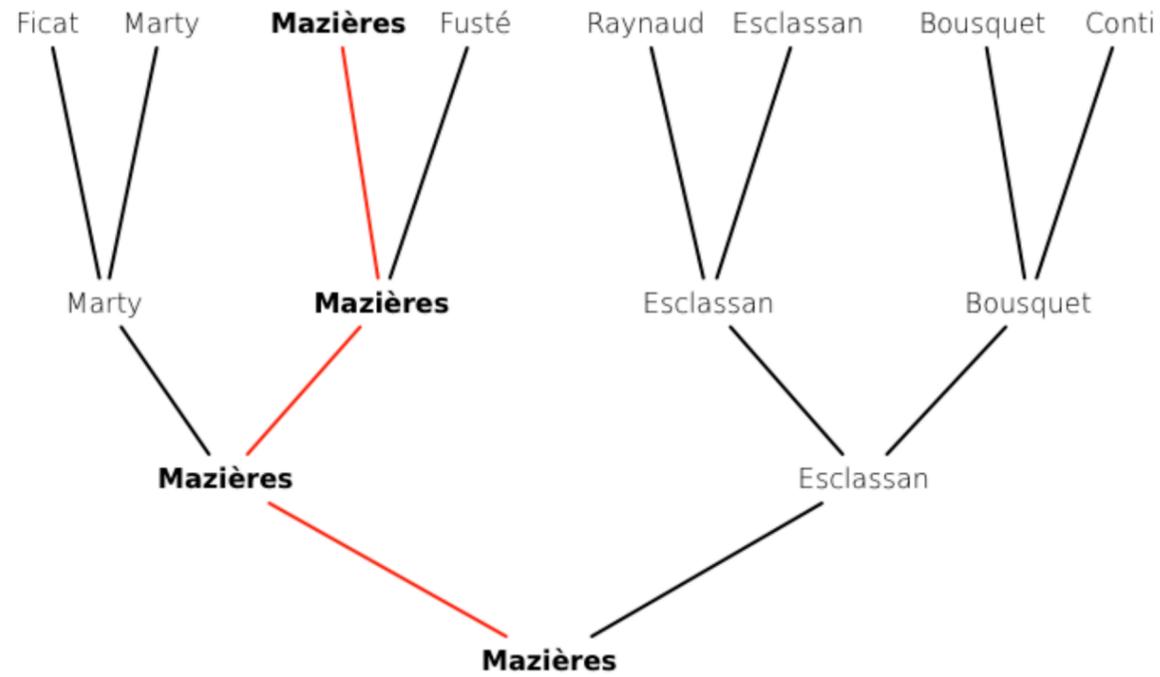
Le nom de famille

Un témoin imparfait du passé

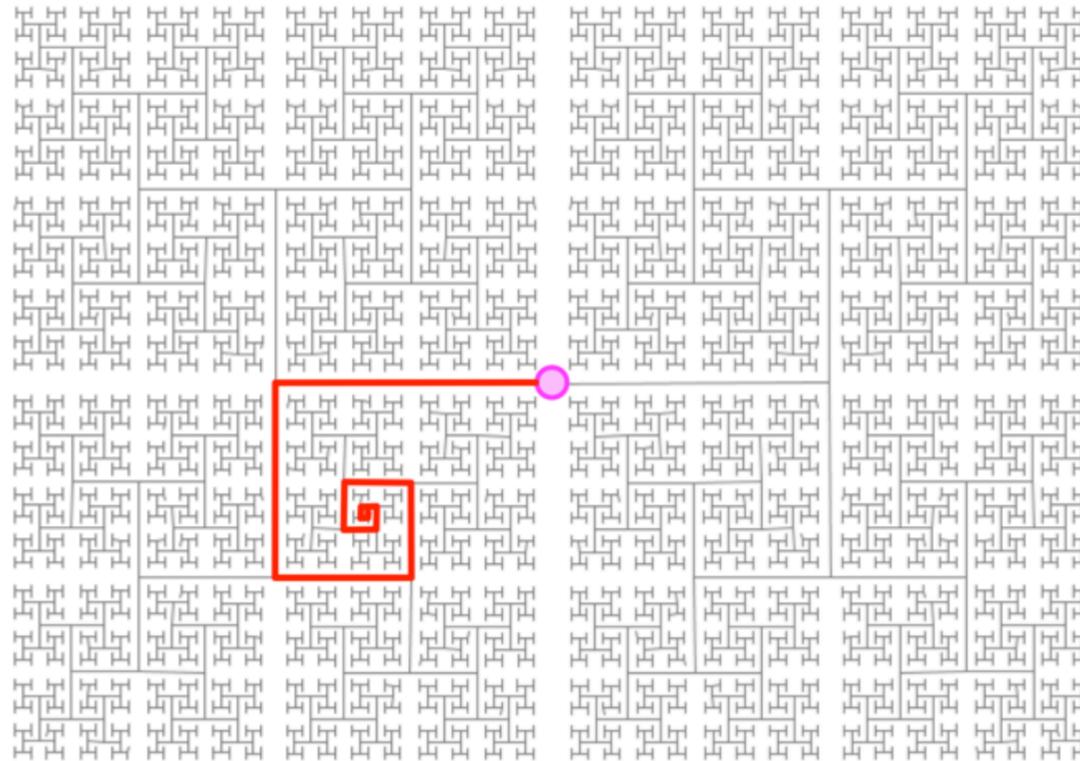
- Outil de distinction : lieux, traits, occupations, surnoms, ...
- En Europe, usage généralisé entre le XIe et XVe siècle.
- Expression d'un rapport de force entre des puissances administratives émergentes.
- Témoin du lieu et de l'époque où ils ont été gelés.
- Nombreux biais : Esclavage, Colonisation, Patriarcat, etc.

Pourquoi les noms ont-ils (encore) un sens ?

4 générations



12 générations



20 générations

> 1 million
d'ancêtres

Pourquoi les noms ont-ils (encore) un sens ?

- Endogamie : Ethnie, religion, génétique, géographique, sociale, économique, etc.
- Usage en démographie, étude des flux migratoires, médecine, marketing, etc.

Construire un modèle de classification des noms

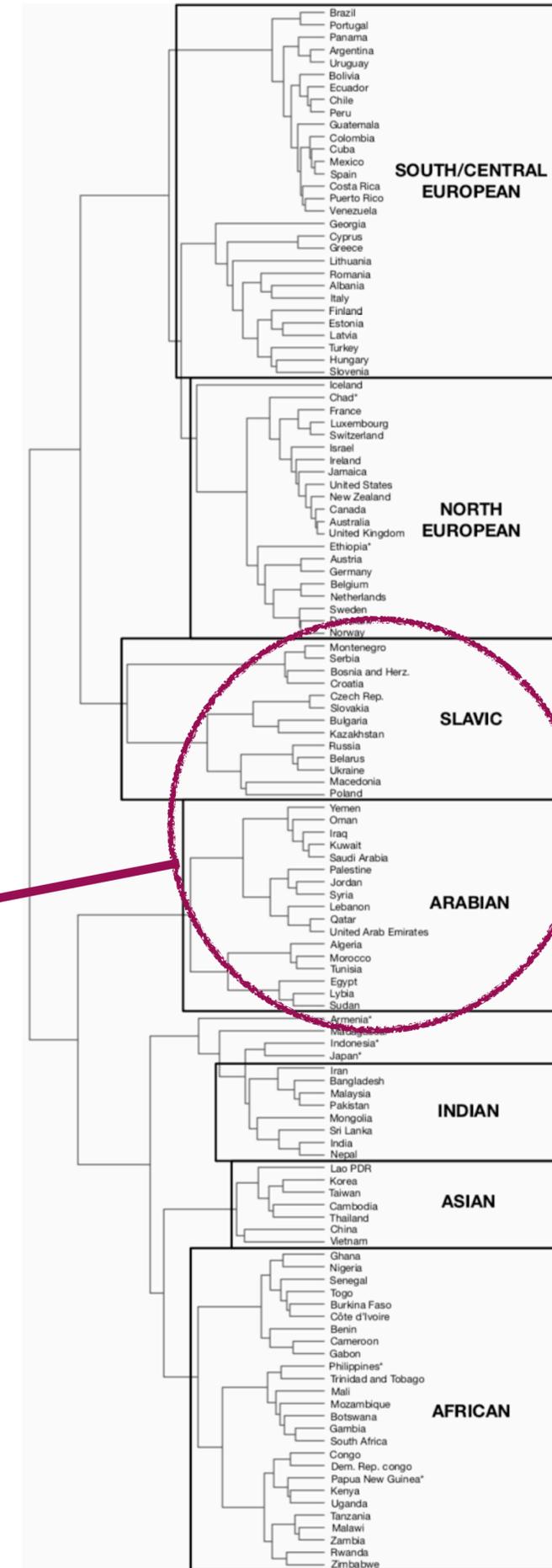
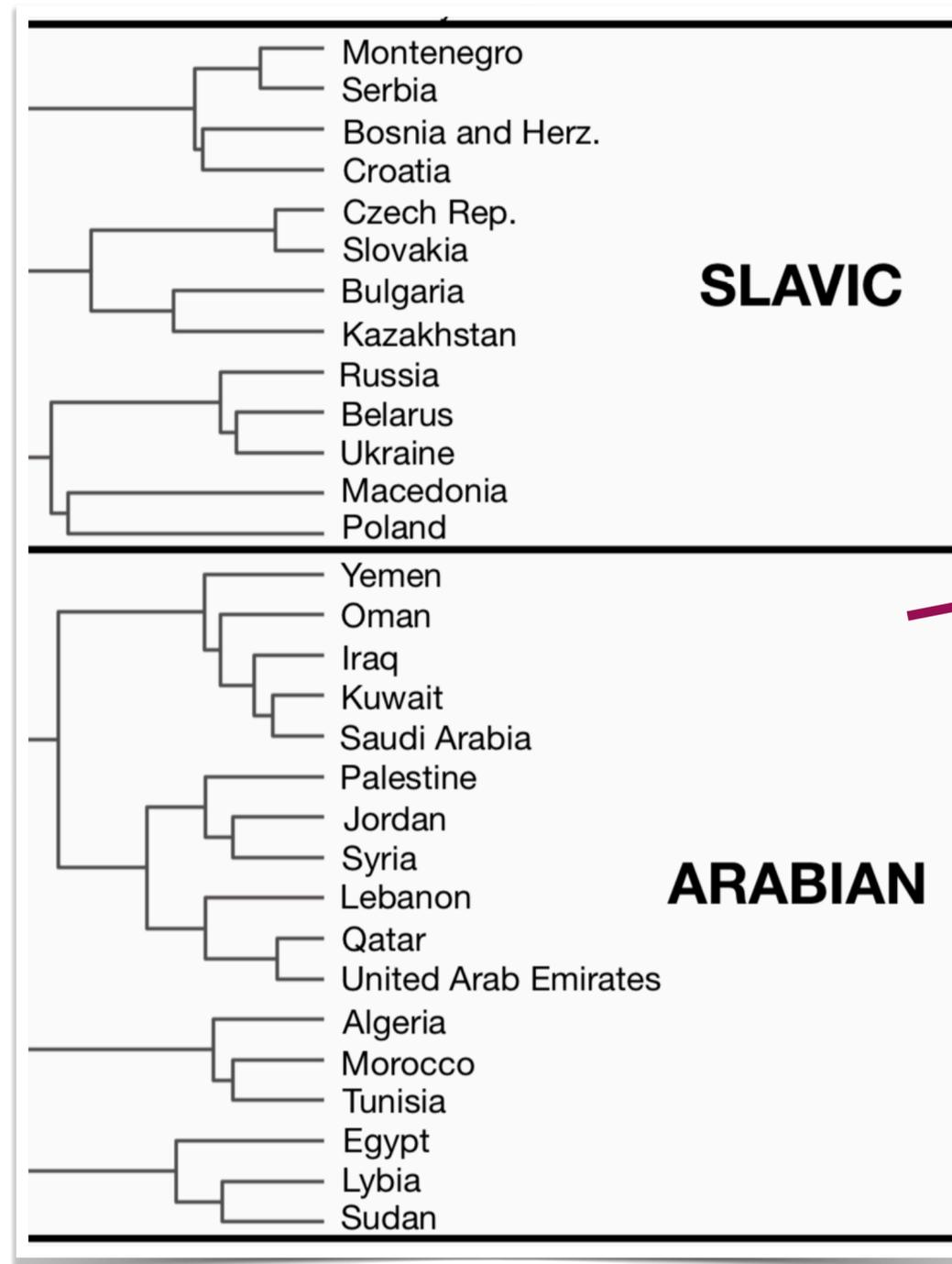
Données et variables

- Base de données : **PubMed**
- **25 millions d'affiliations** : NOM -> PAYS DU LABORATOIRE
- Normalisation par pays, et recherche des noms typiques mais peu “exportés” avec le *Herfindahl-Hirschmann Index* : **650k “core names”**
- Variables avec N-GRAMS : ROTH -> R, O, T, H, ^R, RO, OT, TH, H\$, ^RO, ROT, OTH, TH\$, ...

Construire un modèle de classification des noms

Catégories

- Regroupement des n-grams par pays
- Regroupement hiérarchique (Ward)
- Représentation intuitive des origines avec 7 catégories



Construire un modèle de classification des noms

Apprentissage

- Algorithme : Classification naïve bayésienne.
- Performances hétérogènes :

Cluster	Core names		Class. Perf.	
	<i>Total</i>	<i>Evaluation</i>	<i>Precision</i>	<i>Recall</i>
African	30 748	4 529	0.43	0.61
Arabian	31 272	4 596	0.52	0.72
Asian	44 658	6 754	0.61	0.77
CS-European	189 624	28 668	0.81	0.71
Indian	68 145	10 067	0.63	0.72
N-European	216 465	32 469	0.78	0.62
Slavic	65 259	9 843	0.64	0.84
<i>Total</i>	<i>646 171</i>	<i>96 926</i>		

Applications

15 groupes socio-professionnels en France

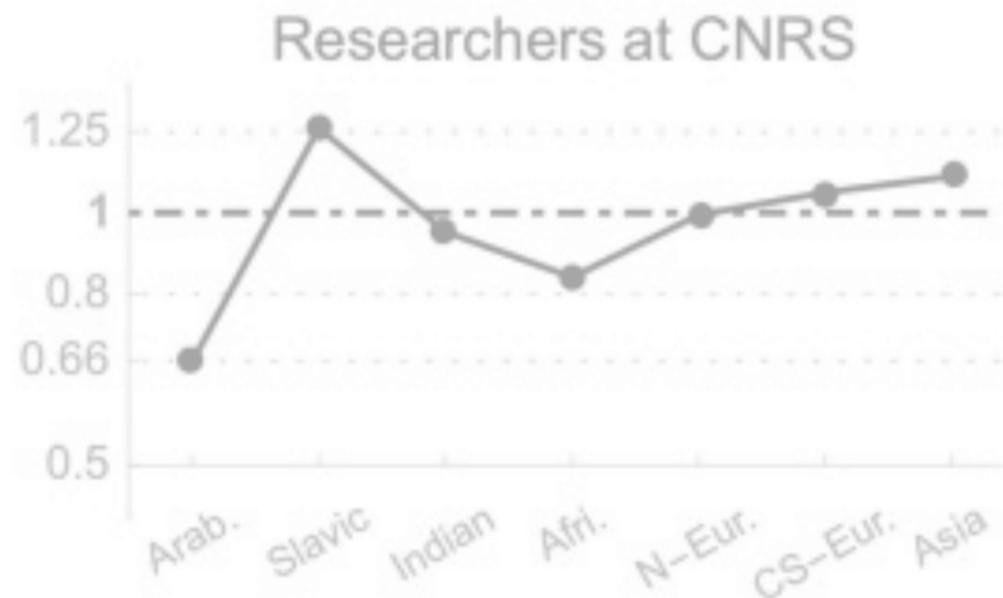
- Quel groupe de référence pour la comparaison ? Brevet
- Comparaison par simple ratio des distributions d'origine
- Cibles :
 - Diplomes (Bac, Bac pro, BTC, CAP, BEP)
 - École polytechniques
 - Corps professionnels (Avocats de Paris, Comptables, Pharmaciens et Vétérinaires)
 - Corps politiques (Députés, Sénateurs, Maires)

Résultats

<https://namograph.antonomase.fr/>

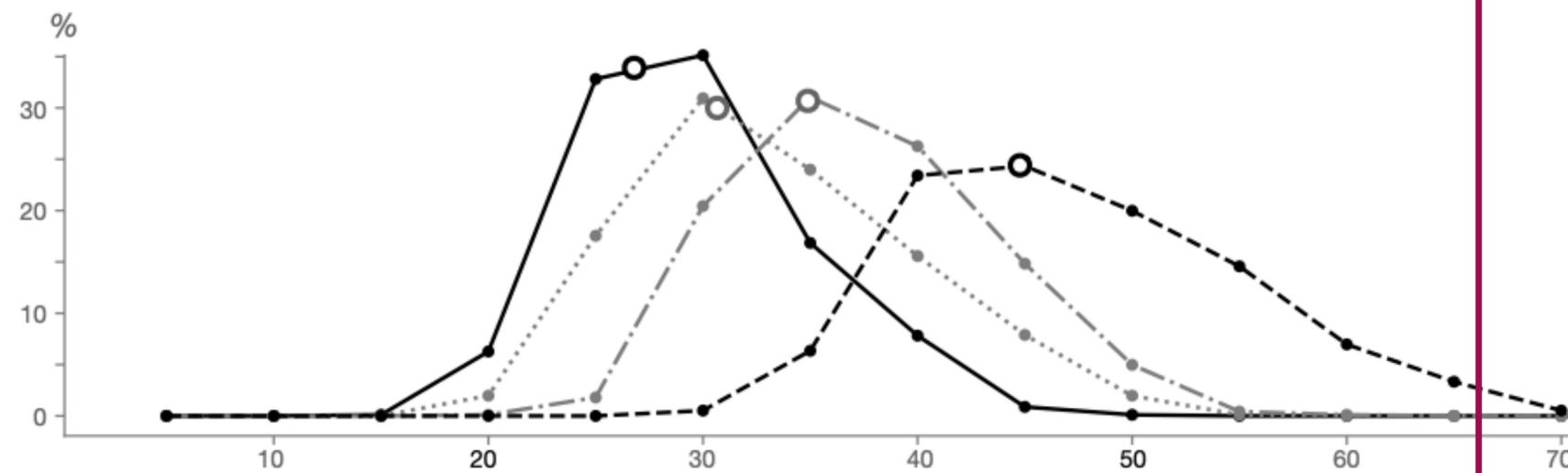
Présentation de deux articles

- Mazieres, Antoine et Roth, Camille (2018). **Large-scale diversity estimation through surname origine inference.** *Bulletin of Sociological Methodology.*



- Mazieres, Antoine, Menezes, Telmo et Roth, Camille (2020). **Computational appraisal of gender representativeness in popular movies.** *Arxiv.*

Evolution of female face ratio in popular movies



Le genre dans les médias de masse

- Longue tradition d'analyse des représentations et de la représentativité de genre depuis les années 40.
- Variables riches
- Faible nombre d'observations
- Méthodologies très changeantes : peu d'études longitudinales

La représentativité de genre dans les films populaires

- Variables simples : Étudier la représentativité et non les représentations
- Utiliser l'analyse d'image (*distant viewing* ?) avec des modèles bien connus : détection de visage et inférence du genre.
- Contenu influent, populaire
- Format constant sur plusieurs décennies.

Définir la base de données

- (L'état de l'art utilise souvent le Box Office (\$) pour définir la popularité)
- Films mis à disposition par YIFY (communauté **Torrent**) avec au moins 3 *seeders* : ~13k films
- ... Et avec un certain nombre de métadonnées sur **IMDb** : Budget (med: 23m\$), Box Office (med: 43m\$), note utilisateur, niveau de censure, durée, genre (romance, crime, ...), année de sortie, etc.
- Total : **3776 films**
- **1985-2019**, minimum 50 films par an.

Application des modèles de détection

- 1 image / 2 secondes : ~12,4m images.
- Modèles de Mathematica 12.
- 6,6m d'images avec visages détectés (~2,5k par film)

Évaluation des inférences



Does the **face** appearing inside the red frame **represent a woman or a man**?

A woman A man I don't know There is no face inside the red frame !

Do you see **other faces** that are clearly identifiable outside the red frame?

Yes No I don't know

(a) Face detection

		Humans	
		<i>Positive</i>	<i>Negative</i>
Model	<i>Positive</i>	977	23
	<i>Negative</i>	137	963

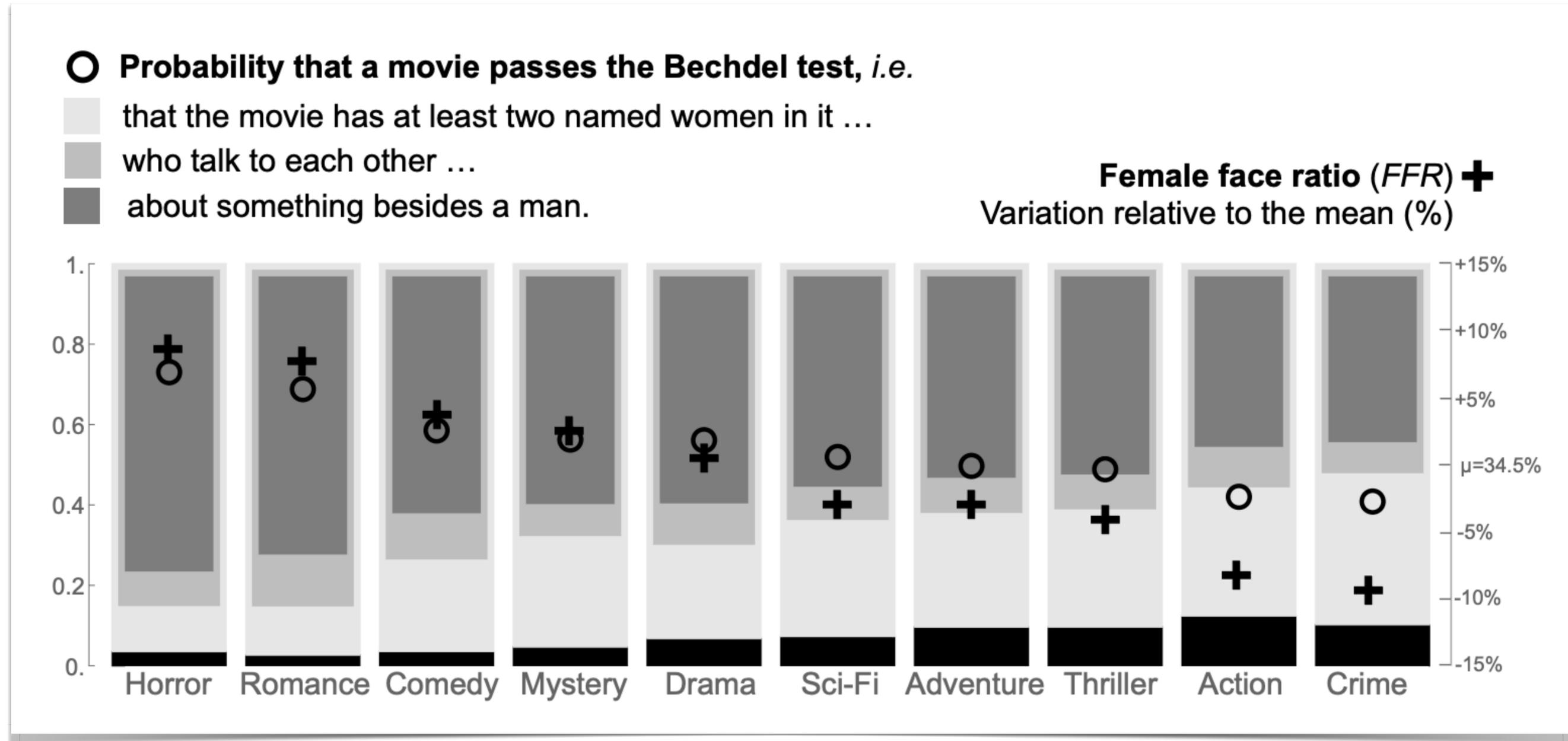
(b) Gender inference

		Humans			
		<i>Female</i>	<i>Male</i>	<i>Doubt</i>	<i>No face</i>
Model	<i>Female</i>	304	162	18	16
	<i>Male</i>	75	410	8	7

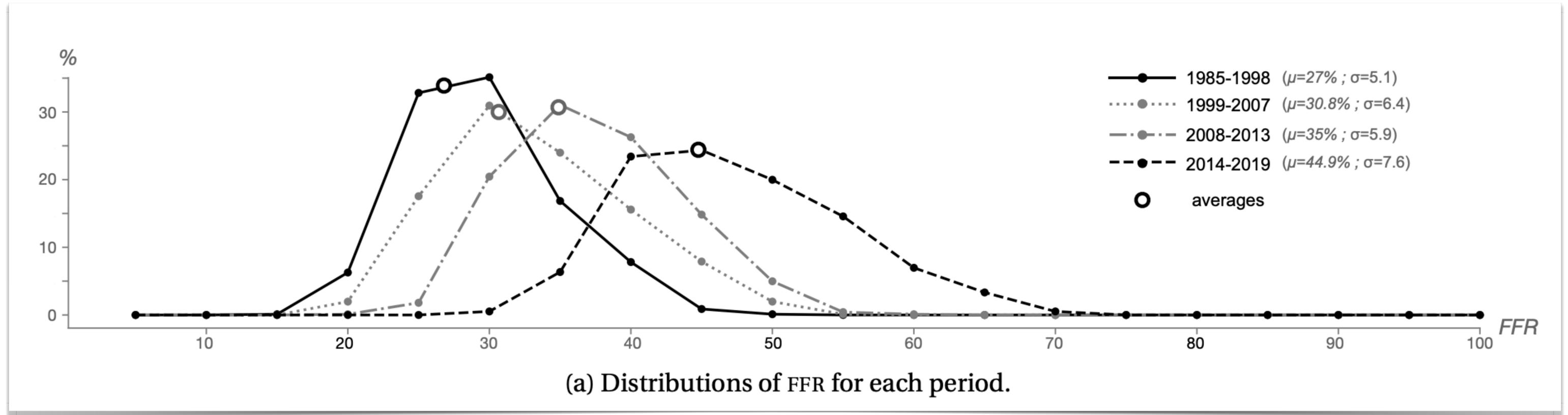
Correction des modèles

- Détection des visages à l'état de l'art (92%)
- Détection de genre en dessous de l'état de l'art (74%)
- Biais de genre : Homme mieux reconnu
- Biais de temps : Images récentes mieux traitées
- Correction bayésiennes pour la mesure du **ratio de visages de femmes (34.52%, $\sigma = 9.19$)**.

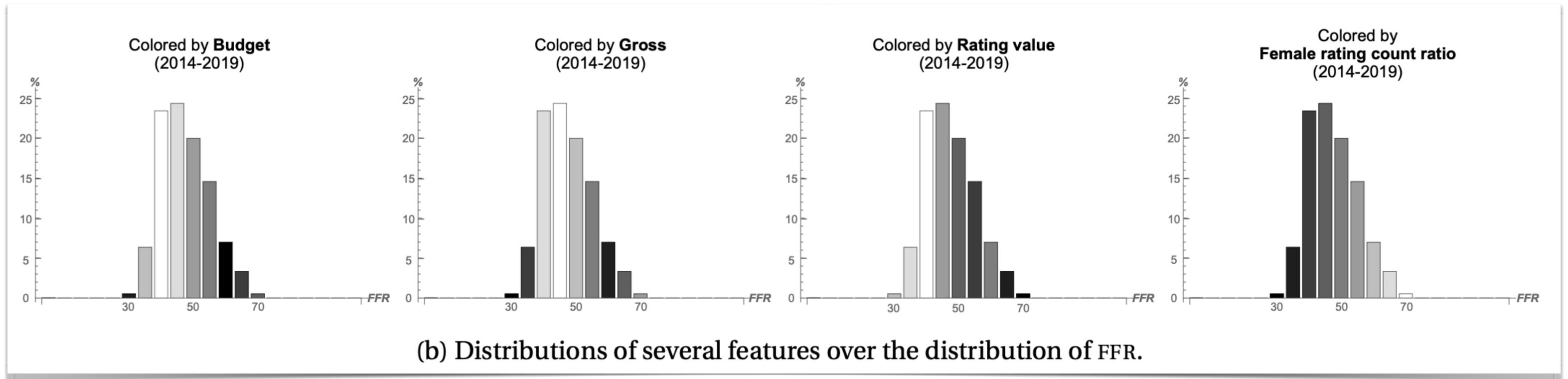
Évaluation plus qualitative



Analyse longitudinale

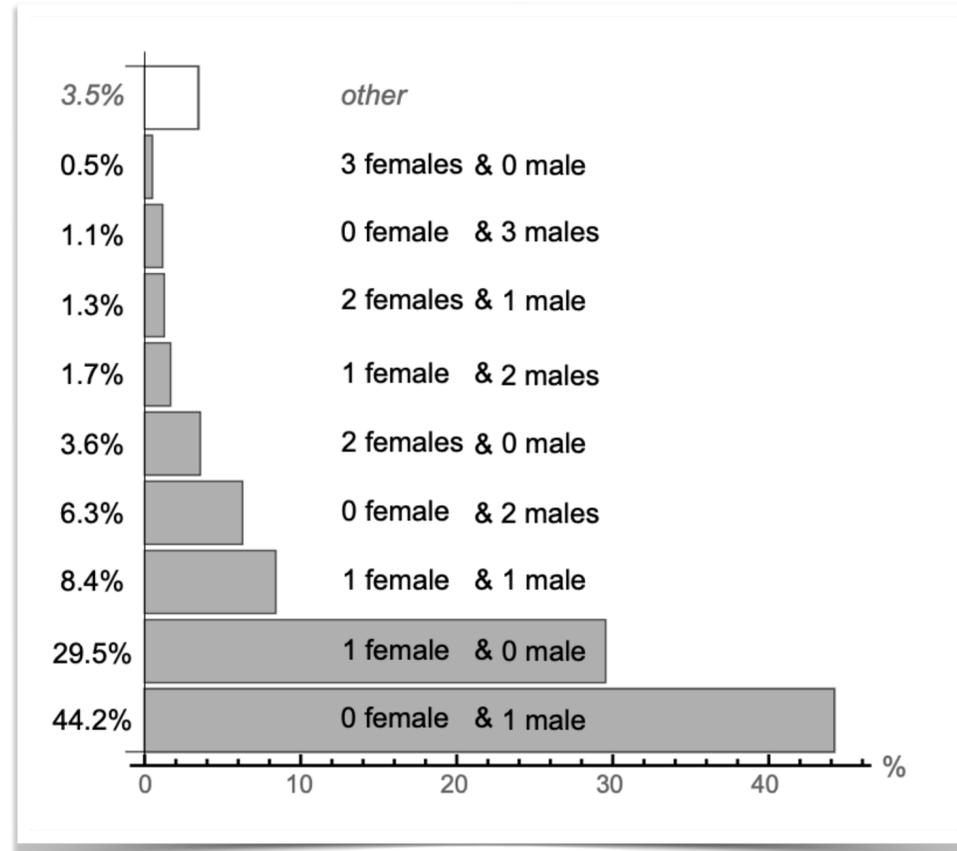


Distributions des métadonnées sur le ratio *FFR*



Mise-en-scène

- Pas de biais perceptible dans la taille des visages (*face-ism*)
- Biais dans la composition des personnages cohérent avec le biais général



Mise-en-cadre

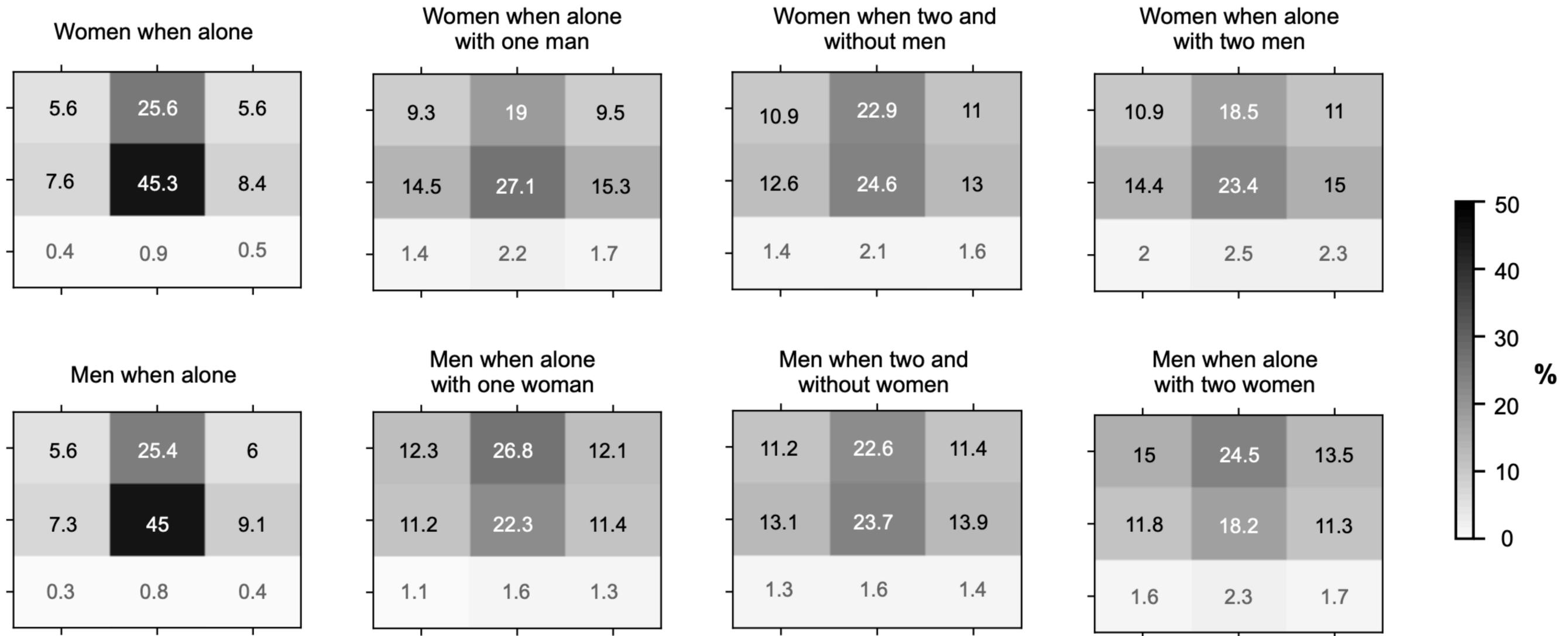


Figure 6: **Distribution of faces position on-screen (2014-2019).**

antoine.mazieres@gmail.com

@mazieres

<https://antonomase.fr>