

Segmentation audio explicable

Distillation de connaissances et factorisation matricielle non-négative

Théo Mariotte, Antonio Almudèvar, Marie Tahon, Alfonso Ortega



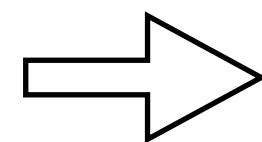
Contexte

Traitement automatique des archives audiovisuelles



Source: <https://arcmc.hypotheses.org/2460>

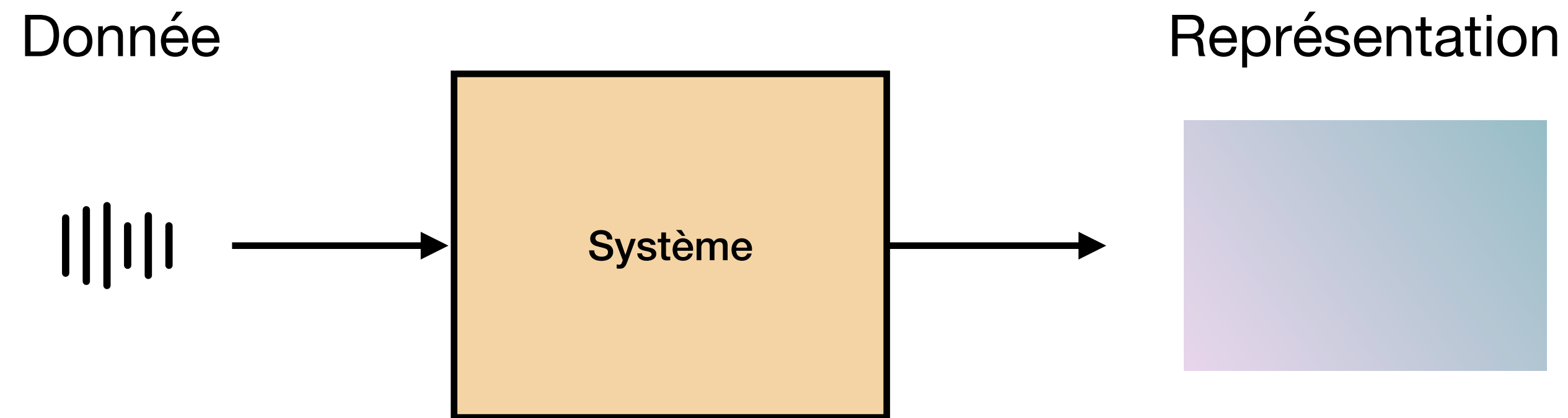
- Grandes bases de données audiovisuelles
- Développement de méthodes de traitement automatique
- Parole : diarization en locuteurs, transcription automatique, reconnaissance du locuteur...
- État de l'art : réseaux de neurones artificiels
- **Problème** : difficile de comprendre l'information utilisée par ces systèmes



Besoin d'expliquer les décisions des systèmes !

Contexte

Expliquer oui... Mais comment ?

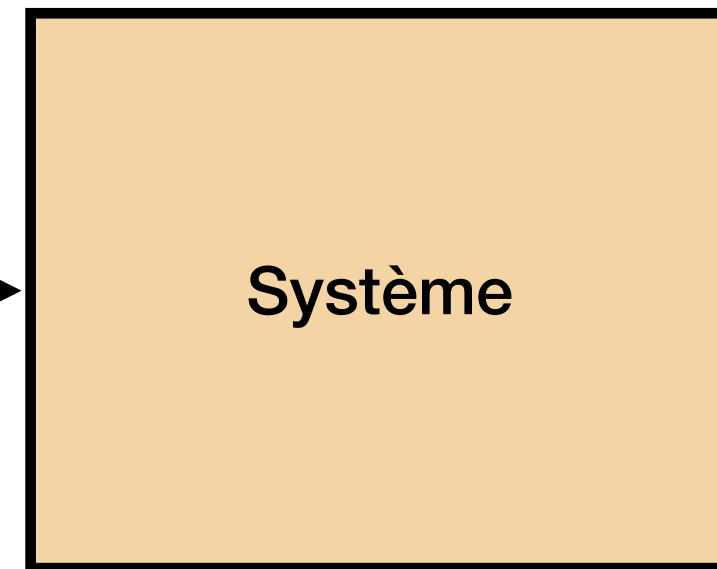
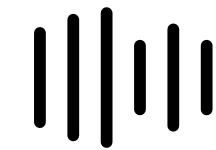


- Relation donnée/représentation abstraite
- Pas de sens « physique »

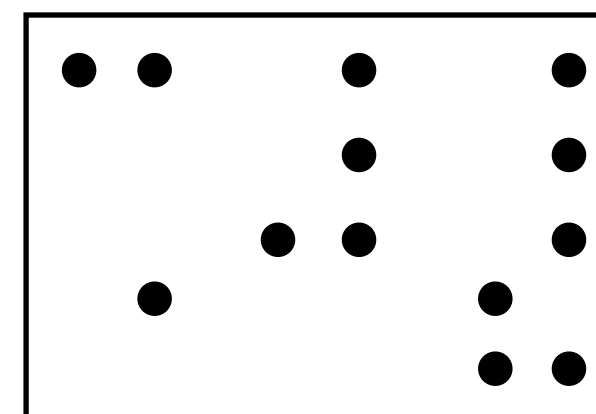
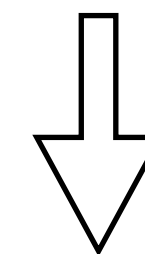
Contexte

Expliquer oui... Mais comment ?

Donnée



Représentation



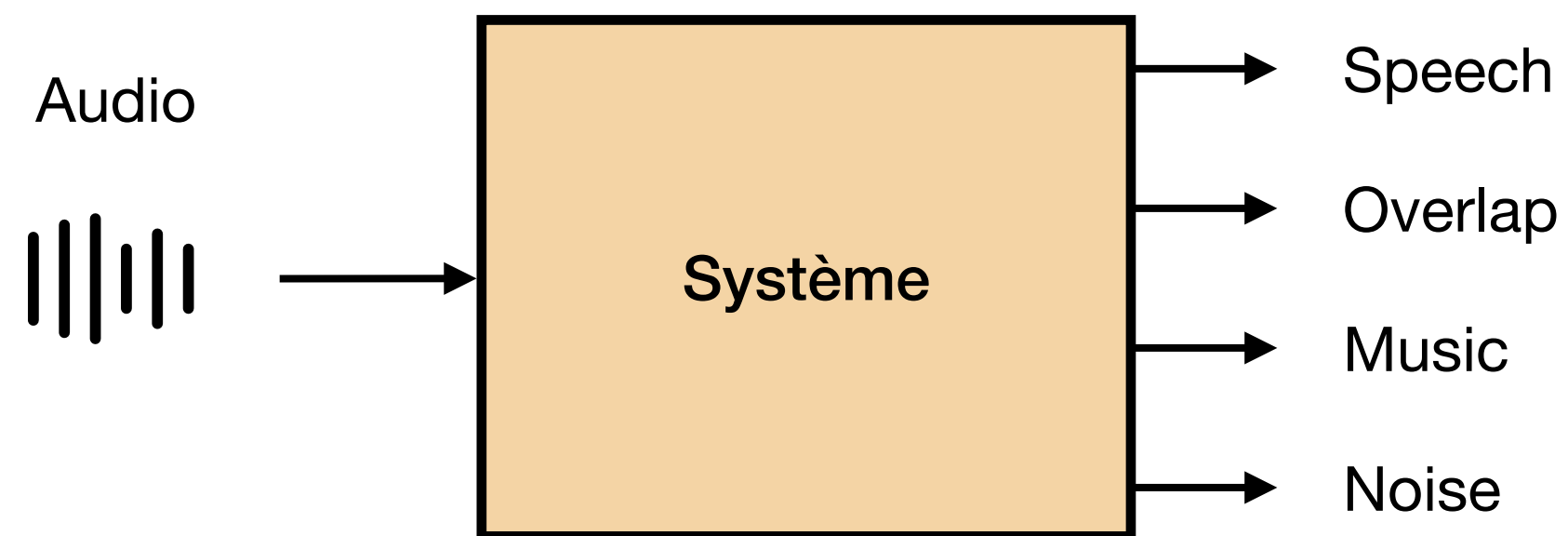
Espace explicable

- Relation donnée/représentation abstraite
- Pas de sens « physique »

- Positivité
- Parcimonie
- Lien vers une représentation « signal »

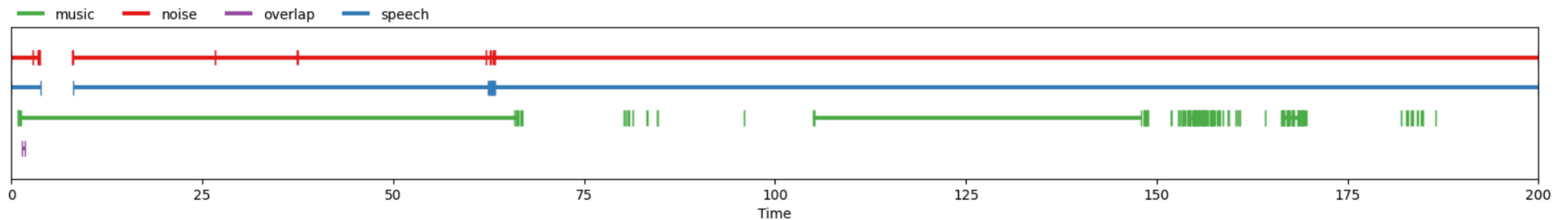
Contexte

Segmentation automatique du signal audio



- Détection de la présence d'évènements en **fonction du temps**
- Utile à de nombreuses tâches de traitement automatique audio
- Évaluation à l'aide du F1-score

Segmentation

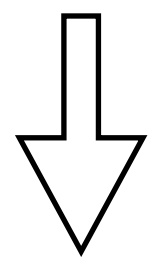
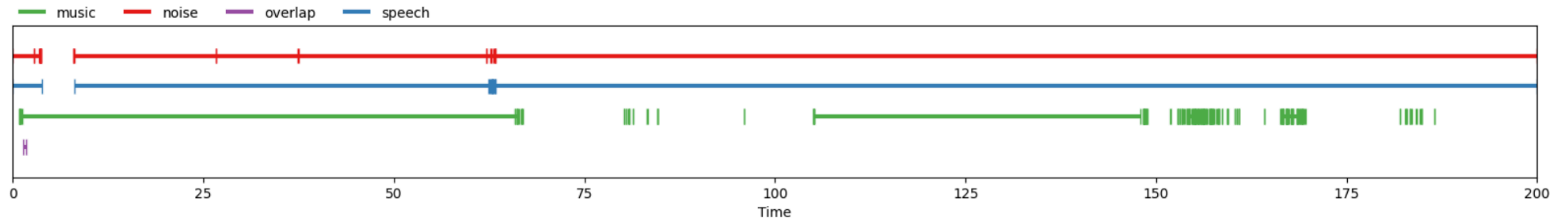


Contexte

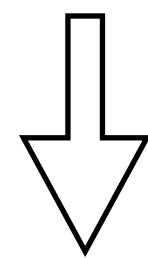
Pourquoi segmenter ?

⇒ Utile à de nombreuses tâches de traitement audio

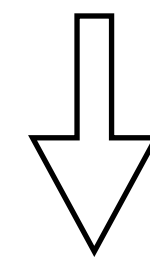
Segmentation



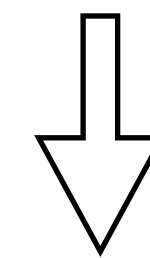
Speaker diarization



Music information retrieval



ASR



Analyse des erreurs

Contexte

Pourquoi expliquer la segmentation ?

- Quelle information est utilisée dans le signal pour segmenter ?
- Comment sont discriminées des classes proches (parole/parole superposée) ?
- Peut on obtenir un prototype de classe ?

Sommaire

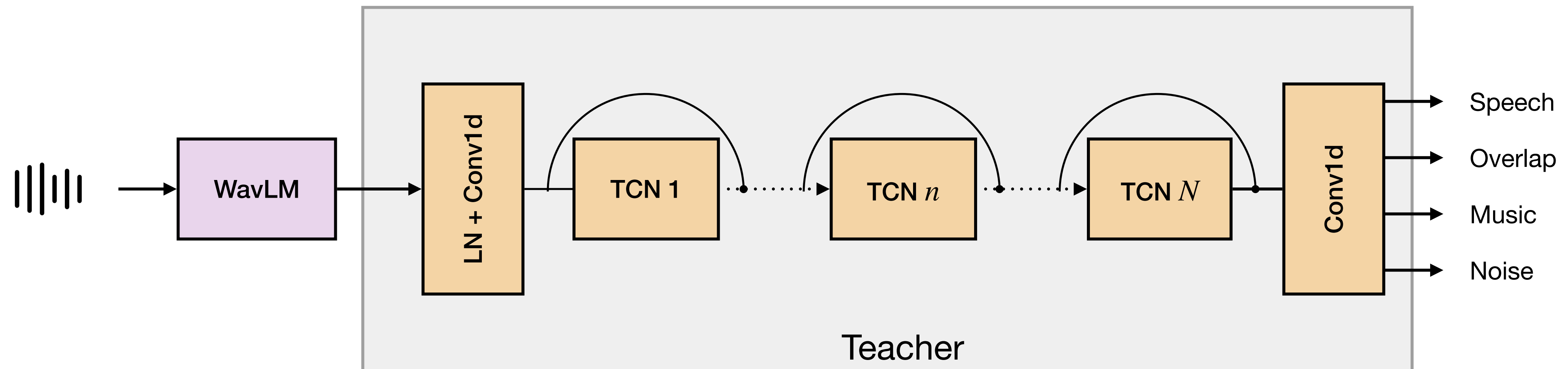
1. Conception d'un système de segmentation audio explicable
 - 1.1. Système *teacher*
 - 1.2. Factorisation matricielle non-négative
 - 1.3. Système *student* explicable
2. Évaluation des performances
 - 2.1. Protocole
 - 2.2. Résultats
3. Extraction d'explications
 - 3.1. Objectif d'explication
 - 3.2. Sélection des composantes pertinentes
 - 3.3. Analyse des explications
4. Conclusions et perspectives

Sommaire

1. Conception d'un système de segmentation audio explicable
 - 1.1. Système *teacher*
 - 1.2. Factorisation matricielle non-négative
 - 1.3. Système *student* explicable
2. Évaluation des performances
 - 2.1. Protocole
 - 2.2. Résultats
3. Extraction d'explications
 - 3.1. Objectif d'explication
 - 3.2. Sélection des composantes pertinentes
 - 3.3. Analyse des explications
4. Conclusions et perspectives

Segmentation audio explicable

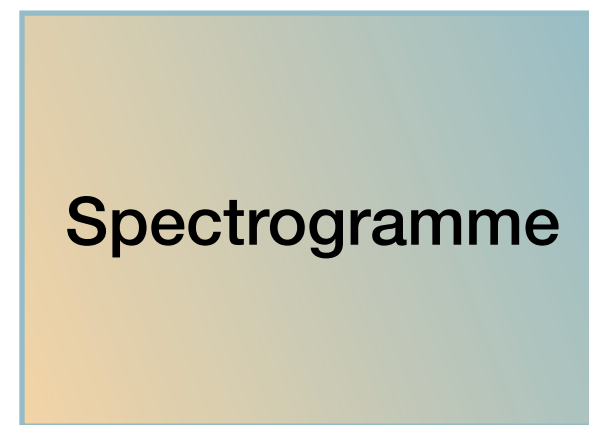
Systeme *teacher*



- Représentation du signal audio avec WavLM (modèle auto-supervisé)
- Modélisation de séquence à l'aide d'un système TCN
- Prédiction des 4 classes à la trame : une prédiction à chaque instant

Segmentation audio explicable

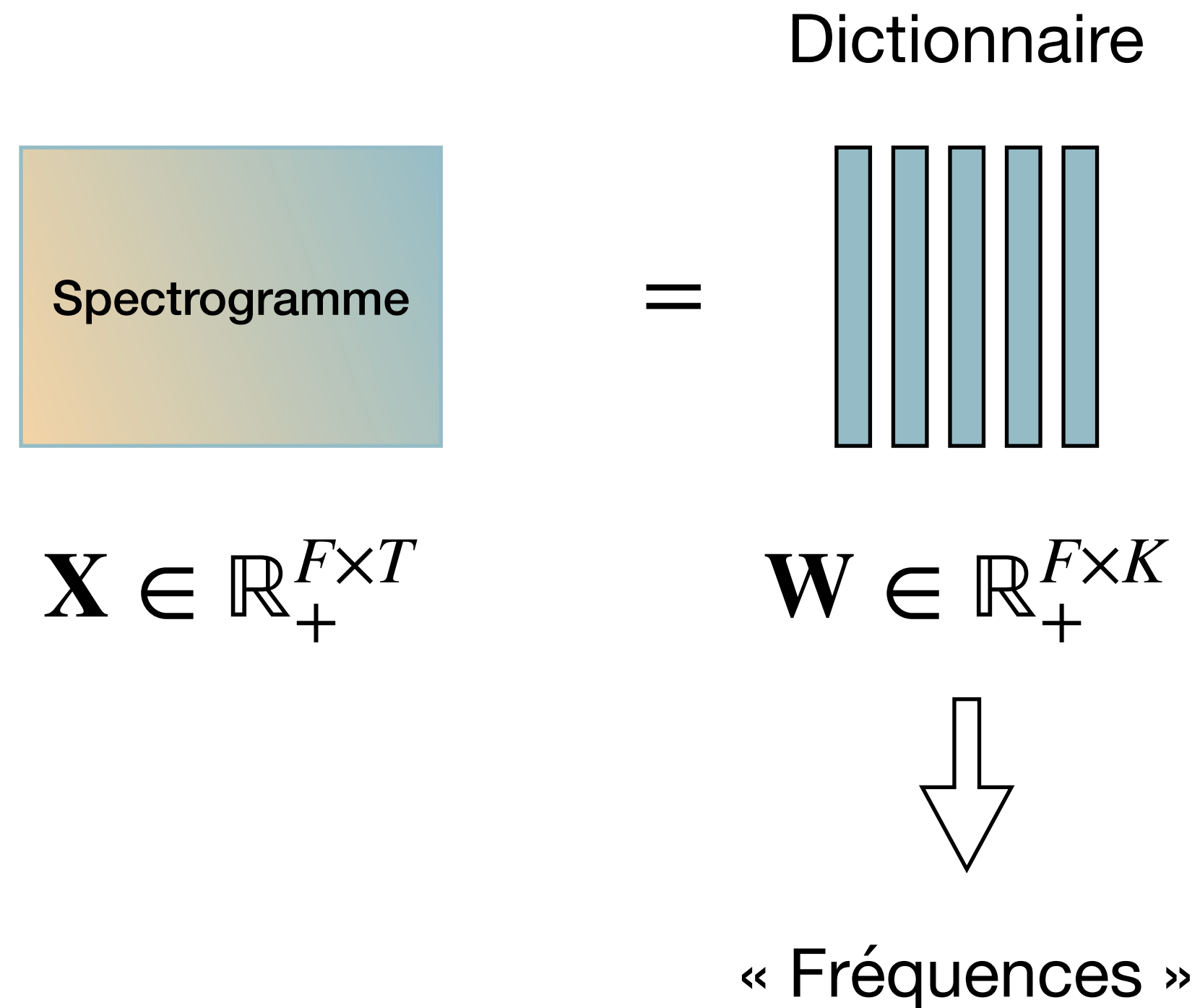
Factorisation matricielle non-négative



$$\mathbf{X} \in \mathbb{R}_+^{F \times T}$$

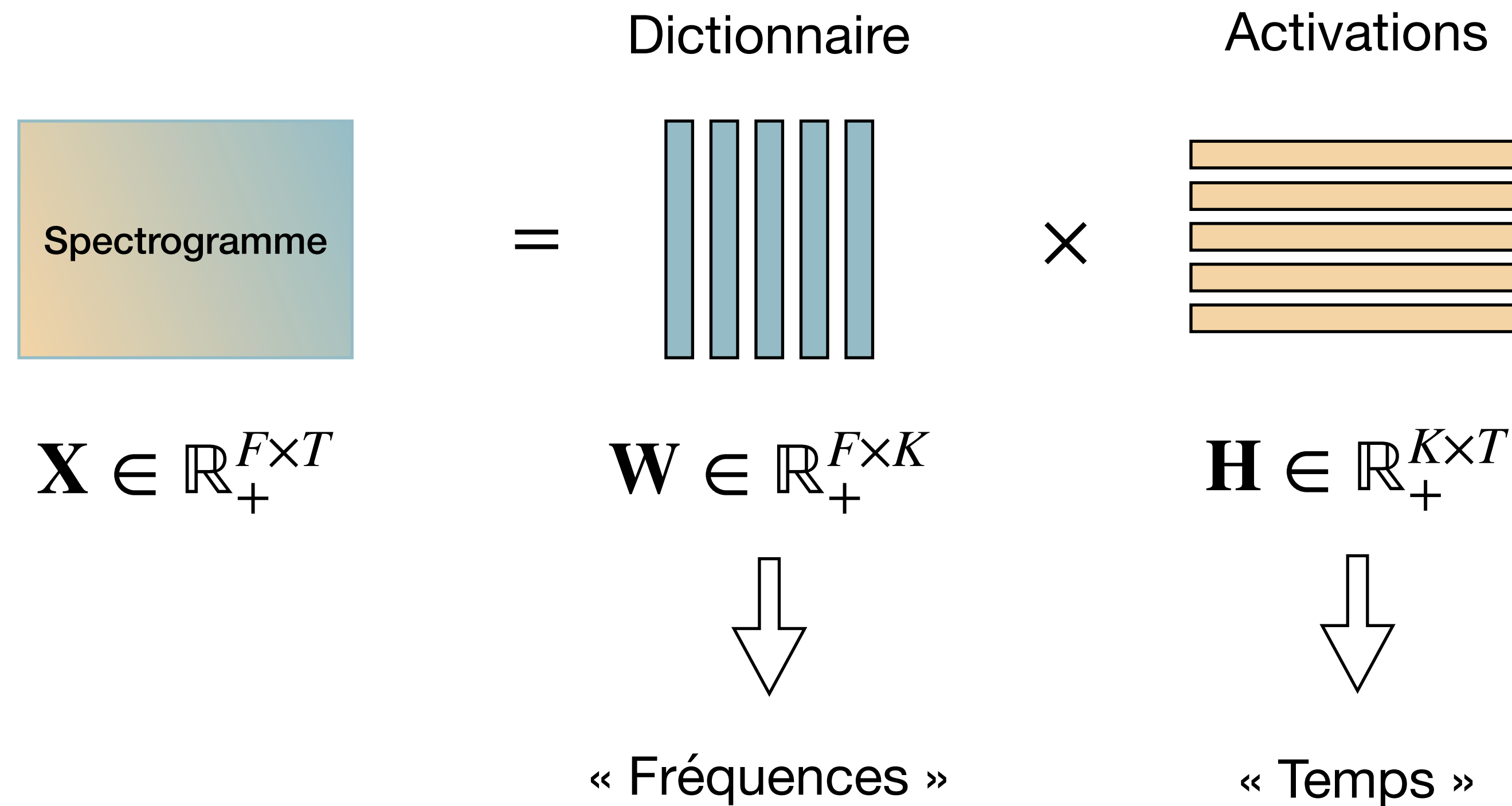
Segmentation audio explicable

Factorisation matricielle non-négative



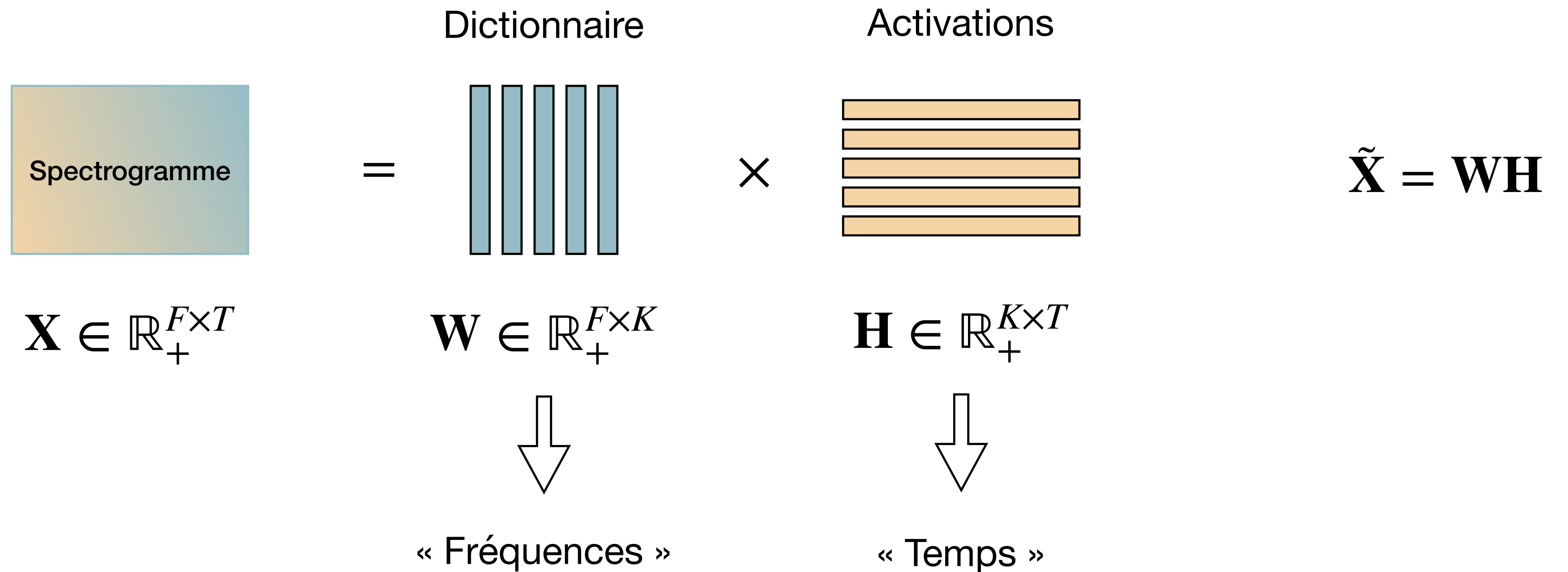
Segmentation audio explicable

Factorisation matricielle non-négative



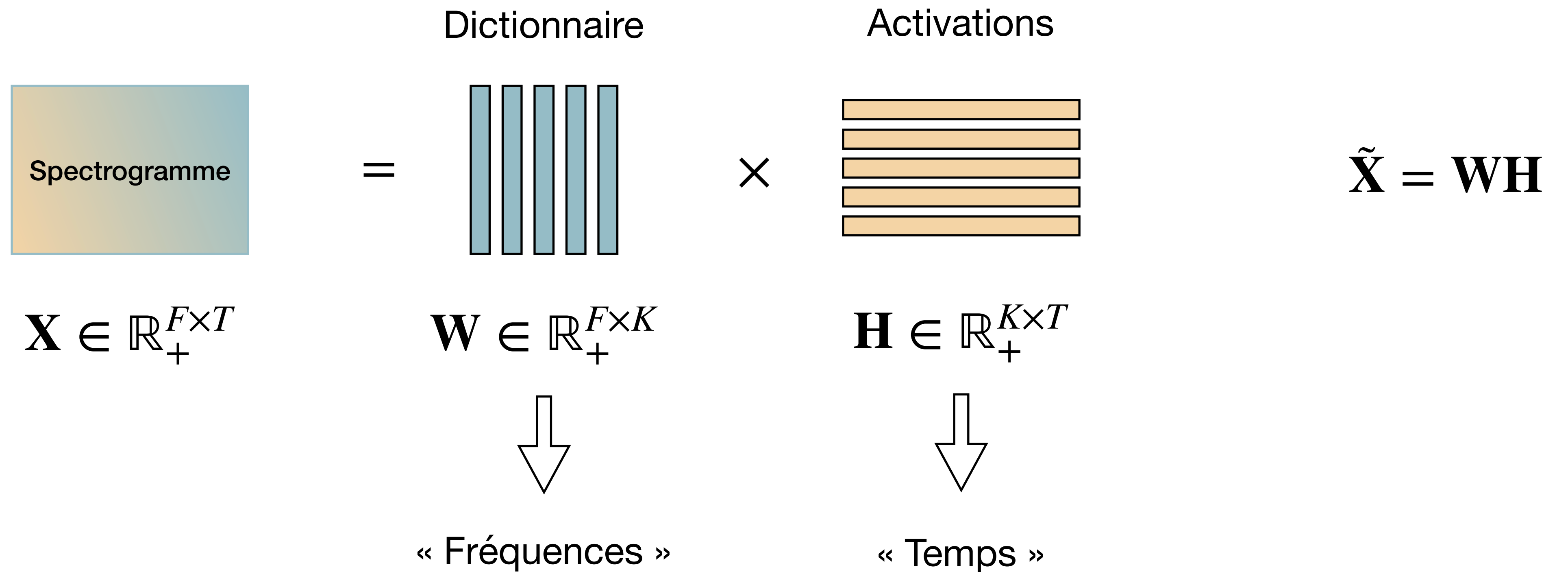
Segmentation audio explicable

Factorisation matricielle non-négative



Segmentation audio explicable

Factorisation matricielle non-négative

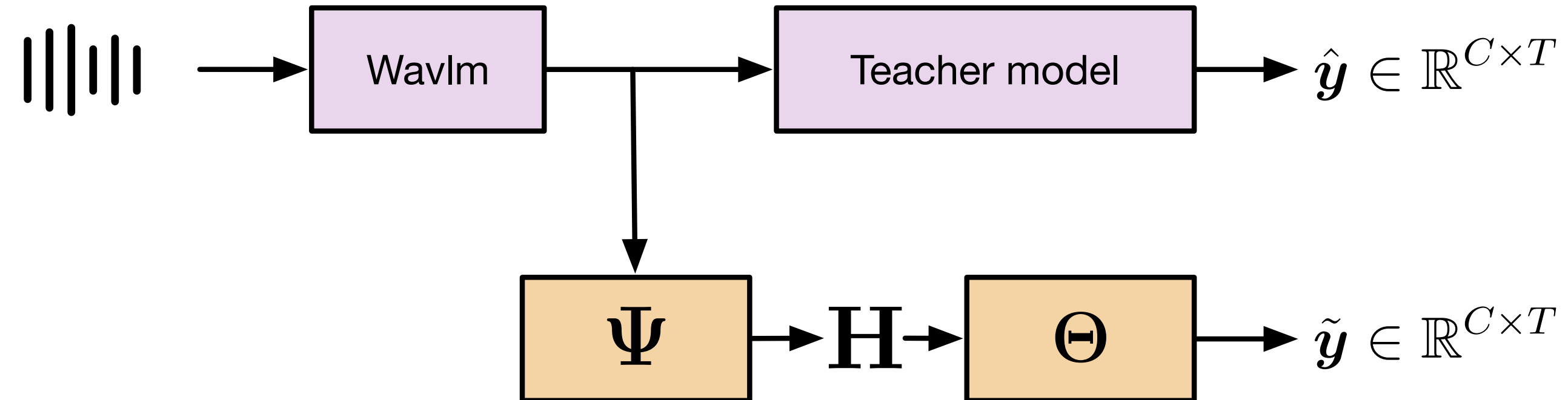


NMF parcimonieuse

$$\min_{\mathbf{W}, \mathbf{H}} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_2 + \mu \|\mathbf{H}\|_1$$

Segmentation audio explicable

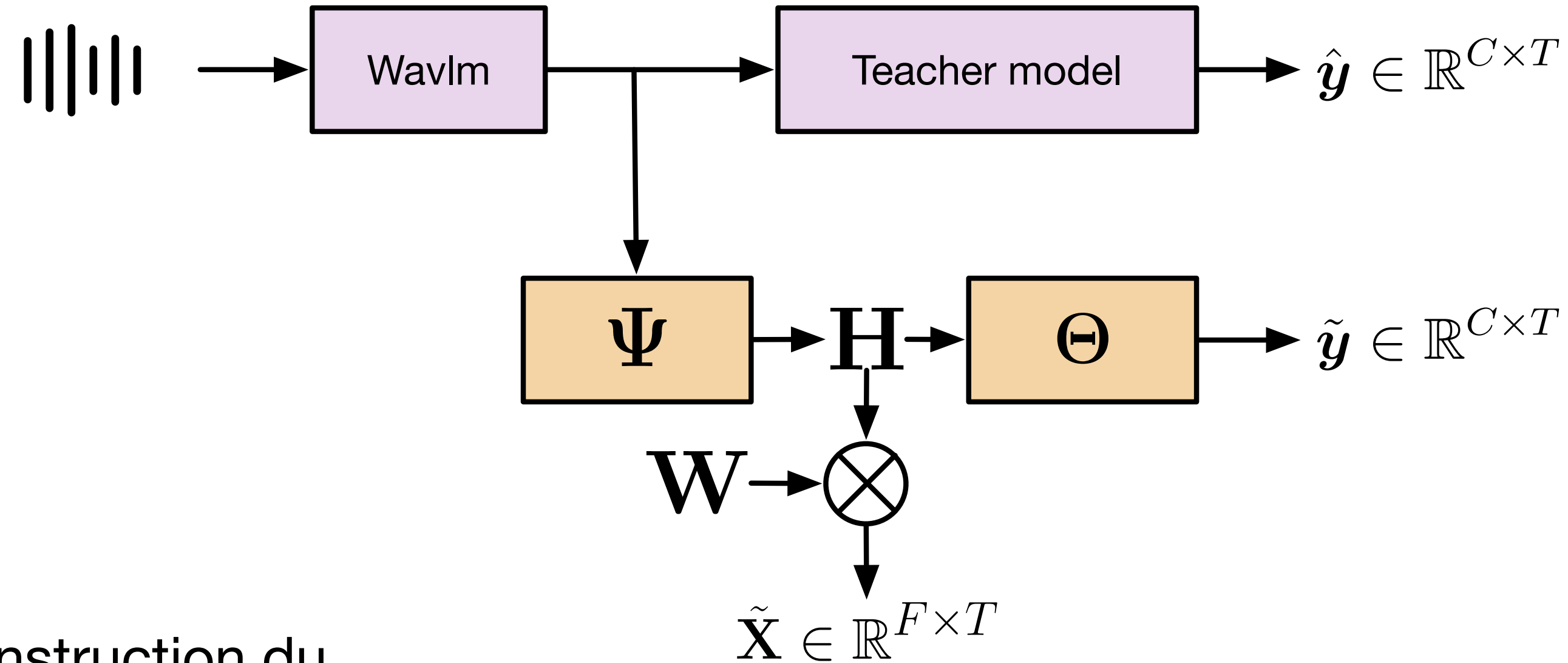
Architecture du *student*



Objectif : contraindre l'espace de représentation \mathbf{H} pour le rendre explicable

Segmentation audio explicable

Architecture du *student*



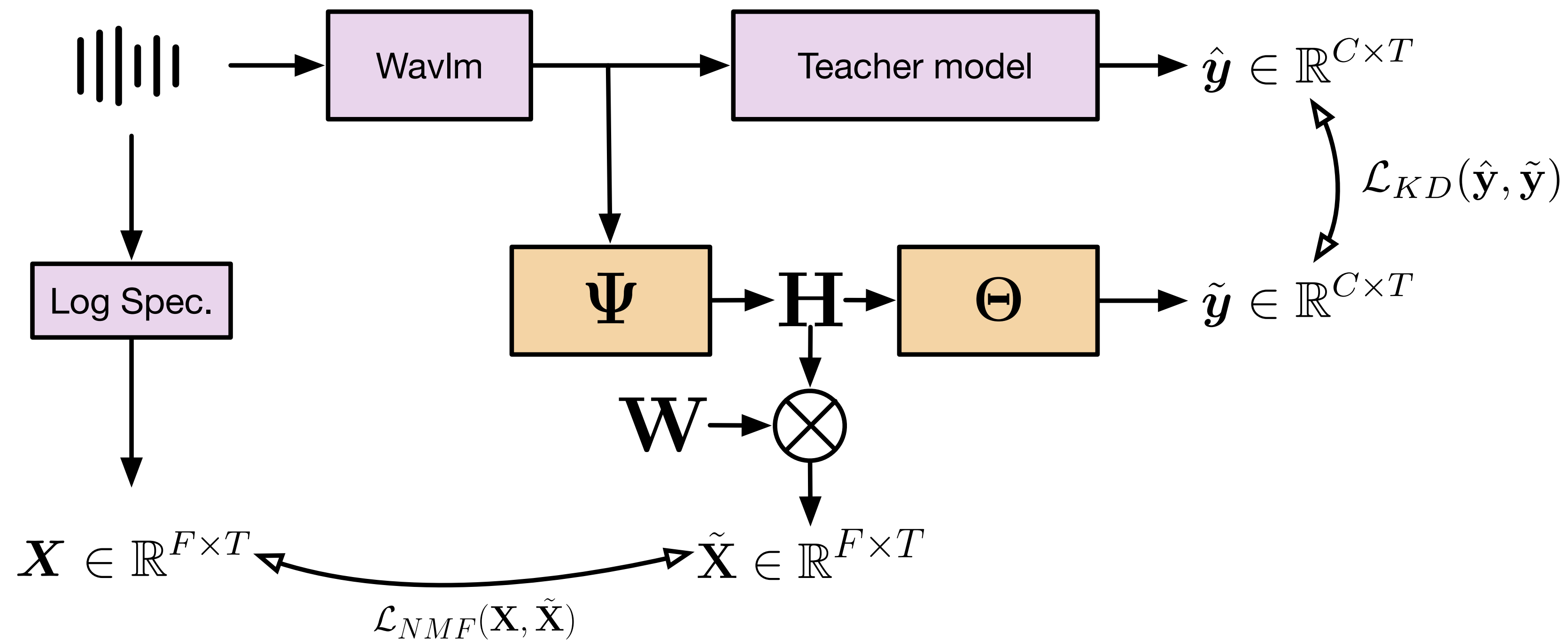
Contrainte 1: reconstruction du spectrogramme

Contrainte 2: non-négativité

Contrainte 3: parcimonie

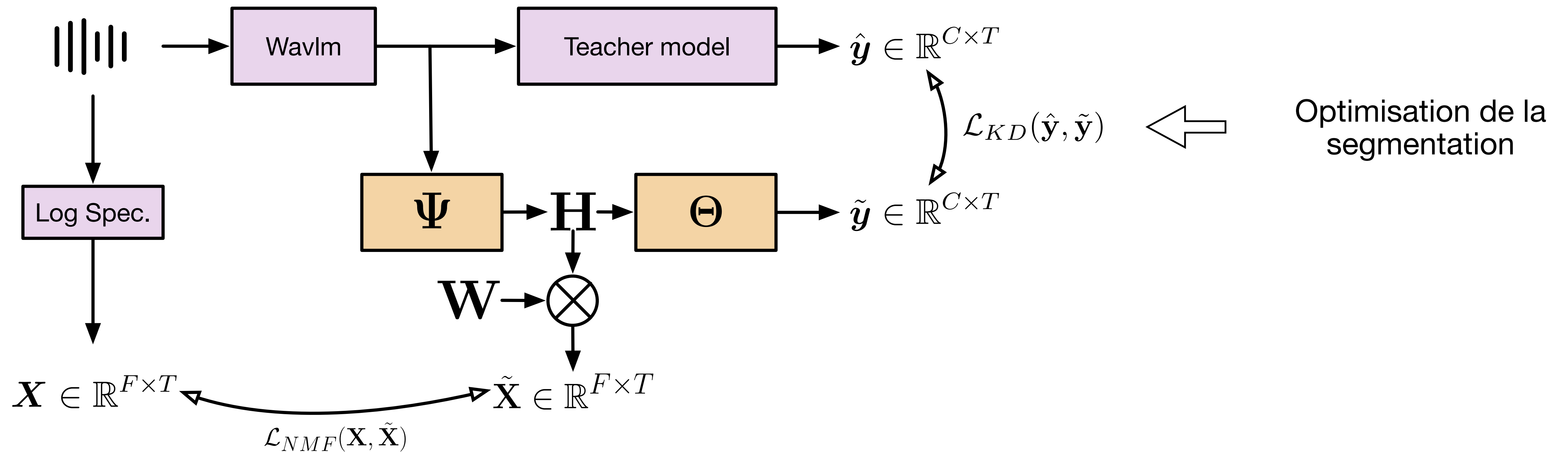
Segmentation audio explicable

Distillation de connaissances



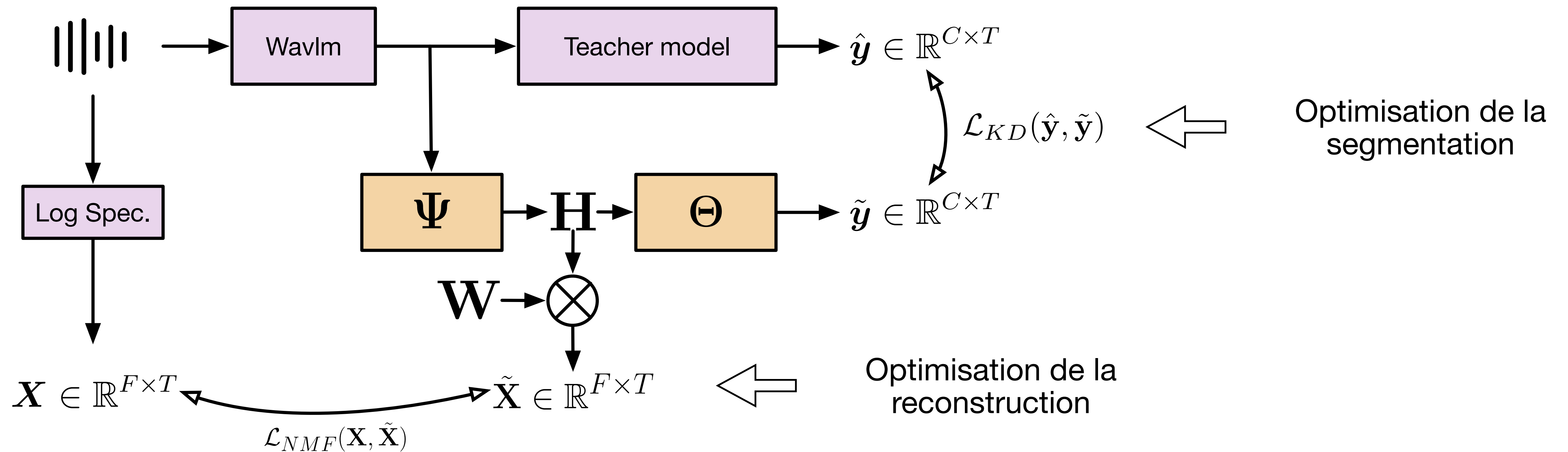
Segmentation audio explicable

Distillation de connaissances



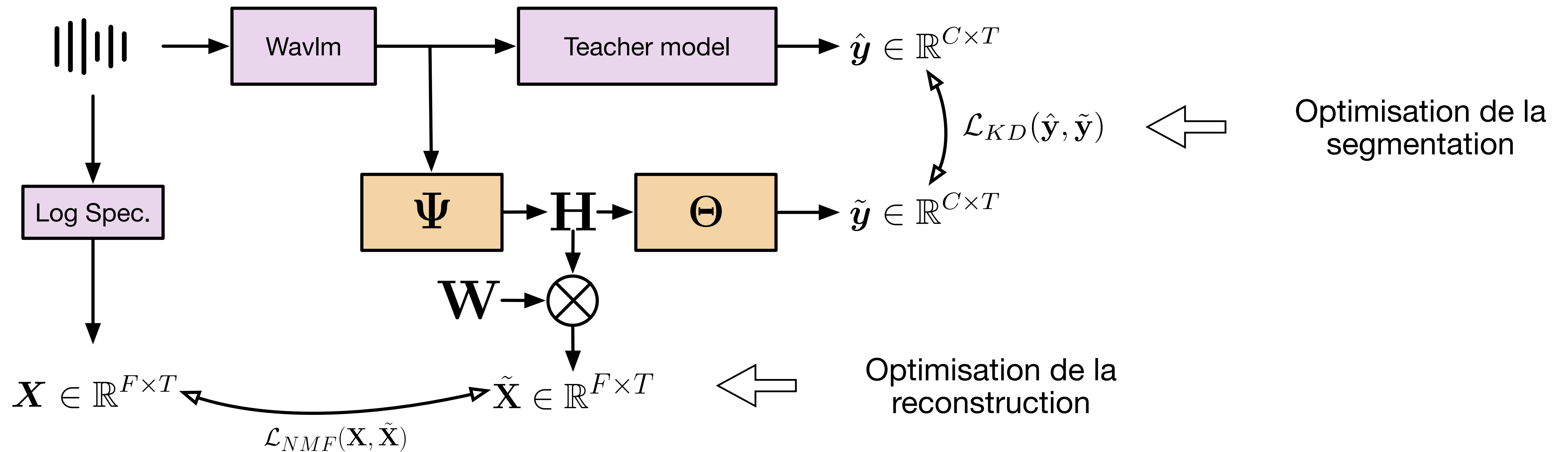
Segmentation audio explicable

Distillation de connaissances



Segmentation audio explicable

Distillation de connaissances



$$\mathcal{L} = \alpha \mathcal{L}_{KD} + \beta \mathcal{L}_{NMF} + \gamma \|\mathbf{H}\|_1$$

Sommaire

1. Conception d'un système de segmentation audio explicable

1.1. Système *teacher*

1.2. Factorisation matricielle non-négative

1.3. Système *student* explicable

2. Évaluation des performances

2.1. Protocole

2.2. Résultats

3. Extraction d'explications

3.1. Objectif d'explication

3.2. Sélection des composantes pertinentes

3.3. Analyse des explications

4. Conclusions et perspectives

Segmentation audio explicable

Protocole d'évaluation

Jeux de données pour l'apprentissage de chaque système :

Dataset	Model		Available label			
	T	P	SAD	MD	ND	OSD
Albayzin	✓	✓	✓	✓	✓	
OpenBMAT	✓			✓		
ALLIES	✓		✓			✓
DiHard III	✓	✓	✓			✓

⇒ Annotations manquantes ignorées pour le calcul de la fonction de perte

⇒ Évaluation des performances sur 3 jeux de données : ALLIES-clean, DIHARD III et Aragon Radio

Segmentation audio explicable

Protocole d'évaluation

Jeux de données pour l'apprentissage de chaque système :

Dataset	Model		Available label			
	T	P	SAD	MD	ND	OSD
Albayzin	✓	✓	✓	✓	✓	
OpenBMAT	✓			✓		
ALLIES	✓		✓			✓
DiHard III	✓	✓	✓			✓

Teacher

⇒ Annotations manquantes ignorées pour le calcul de la fonction de perte

⇒ Évaluation des performances sur 3 jeux de données : ALLIES-clean, DIHARD III et Aragon Radio

Segmentation audio explicable

Protocole d'évaluation

Jeux de données pour l'apprentissage de chaque système :

Dataset	Model		Available label			
	T	P	SAD	MD	ND	OSD
Albayzin	✓	✓	✓	✓	✓	
OpenBMAT	✓			✓		
ALLIES	✓		✓			✓
DiHard III	✓	✓	✓			✓

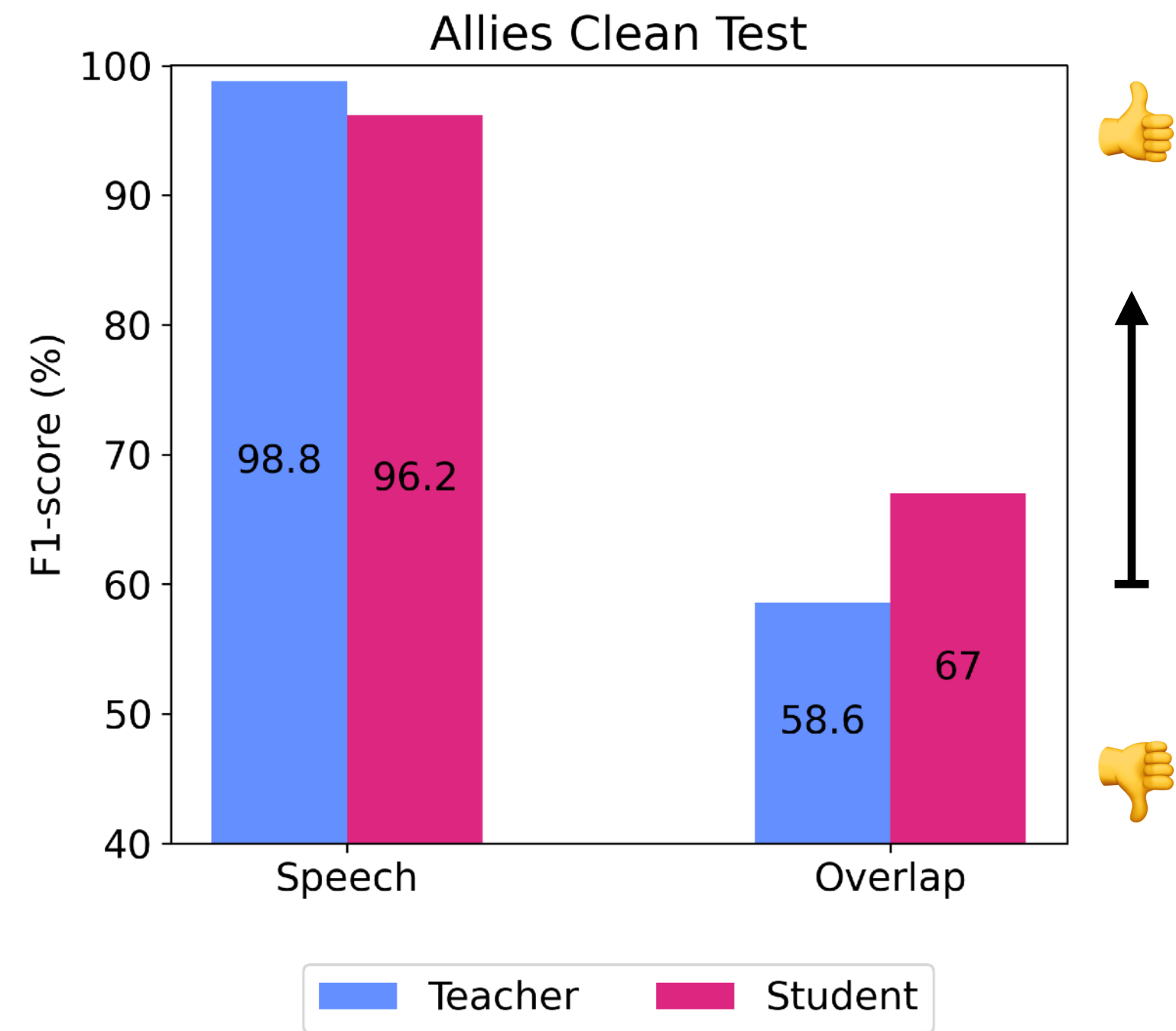
Teacher Student

⇒ Annotations manquantes ignorées pour le calcul de la fonction de perte

⇒ Évaluation des performances sur 3 jeux de données : ALLIES-clean, DIHARD III et Aragon Radio

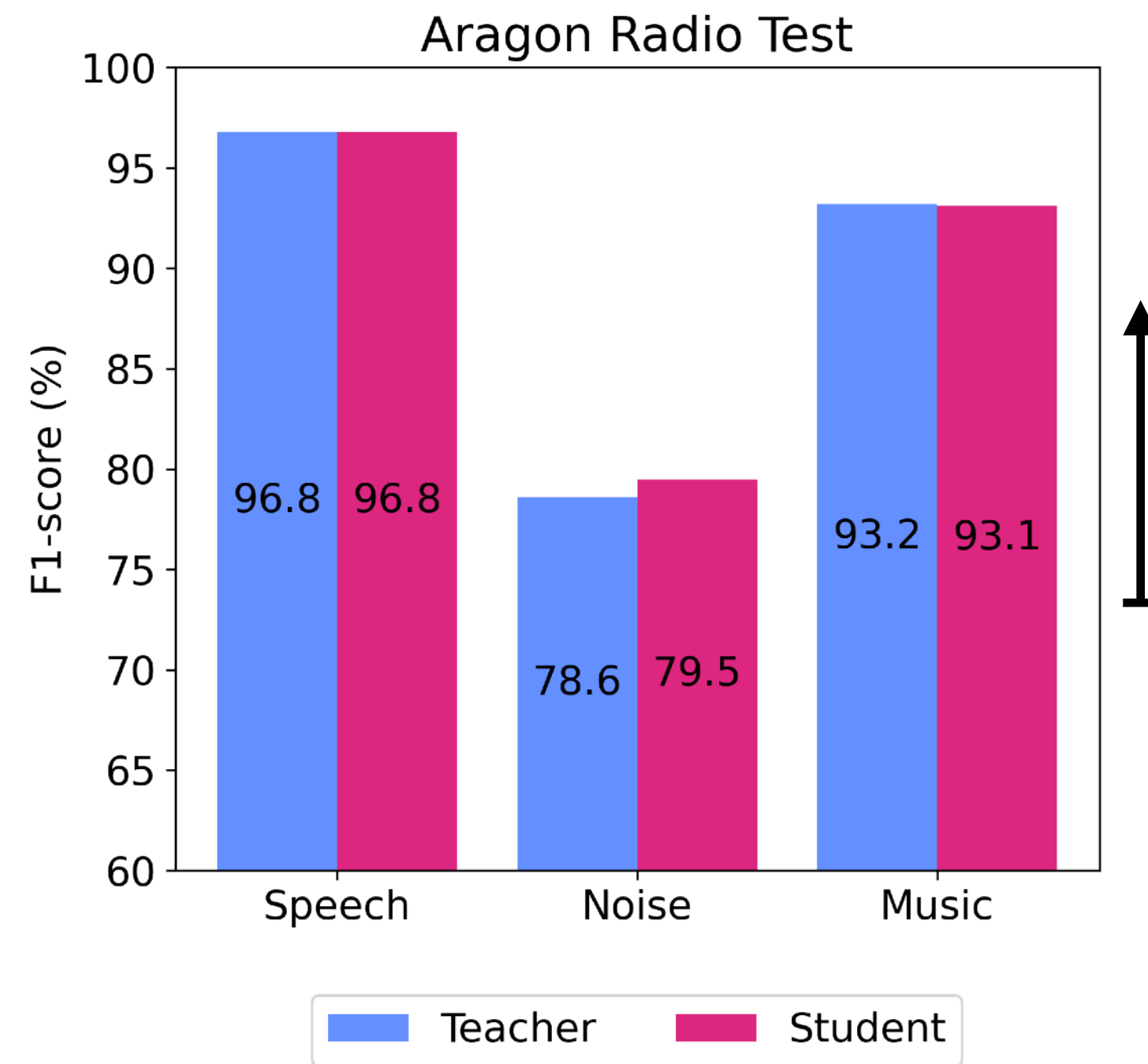
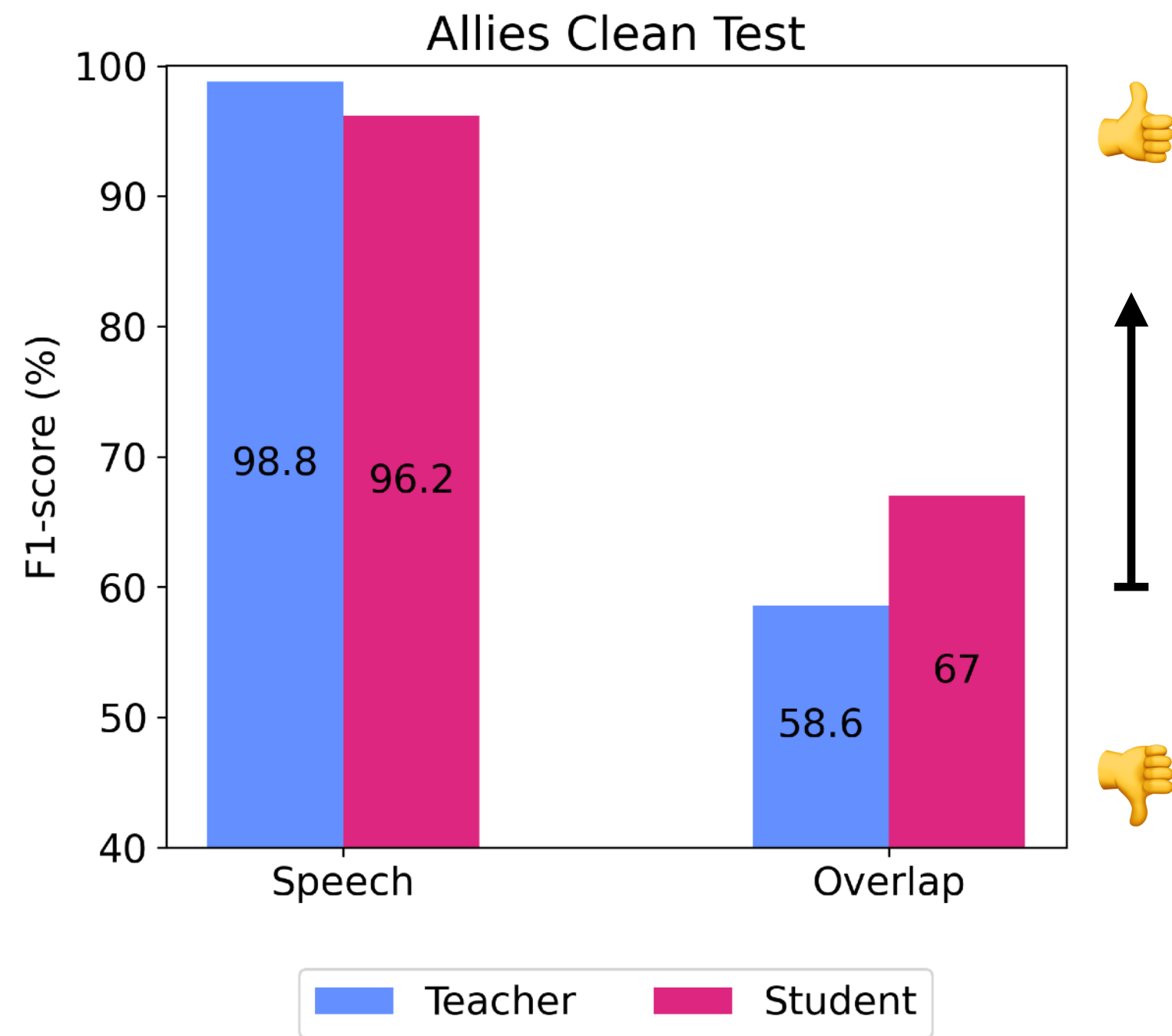
Performances de segmentation

Teacher vs. Student



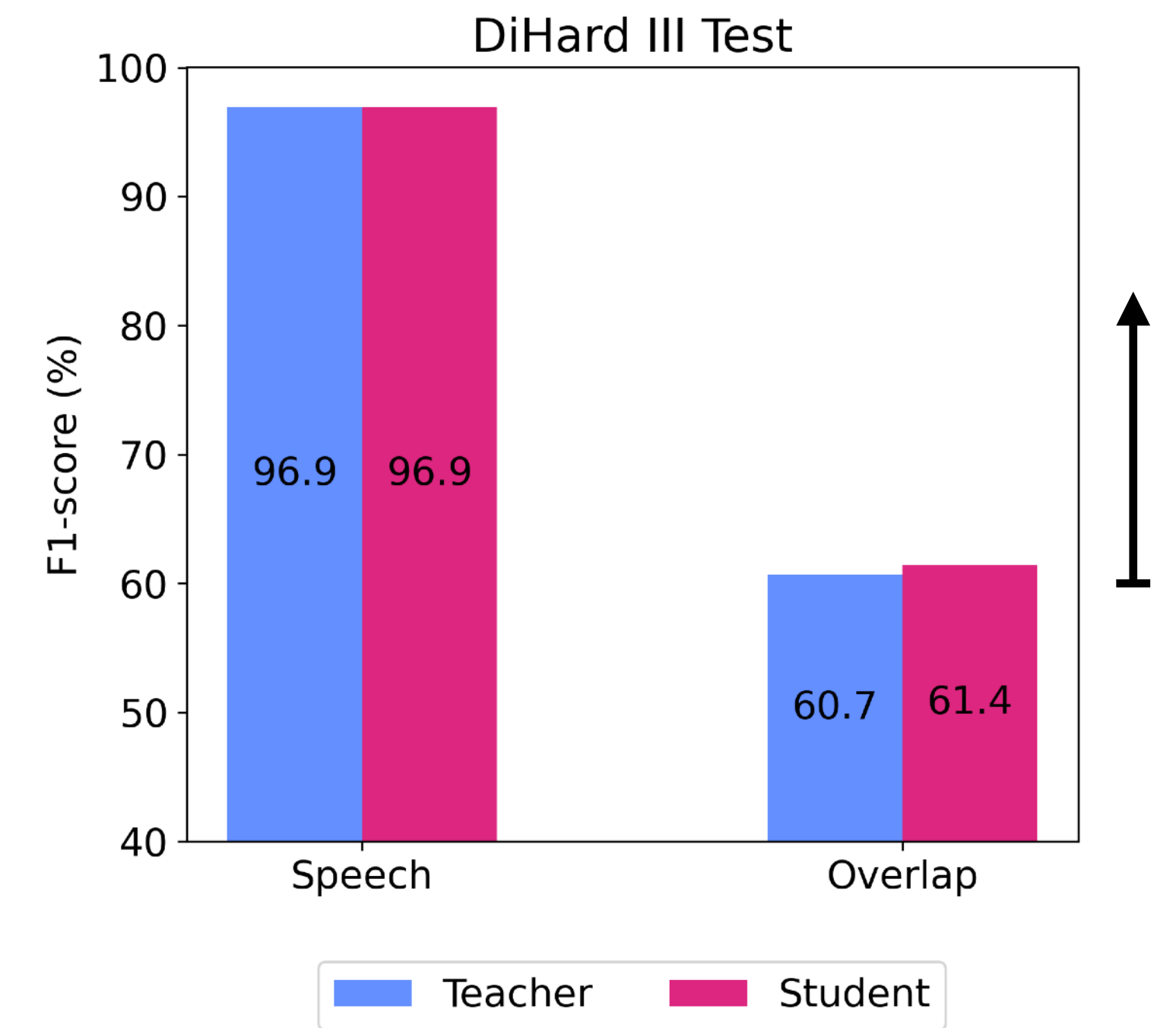
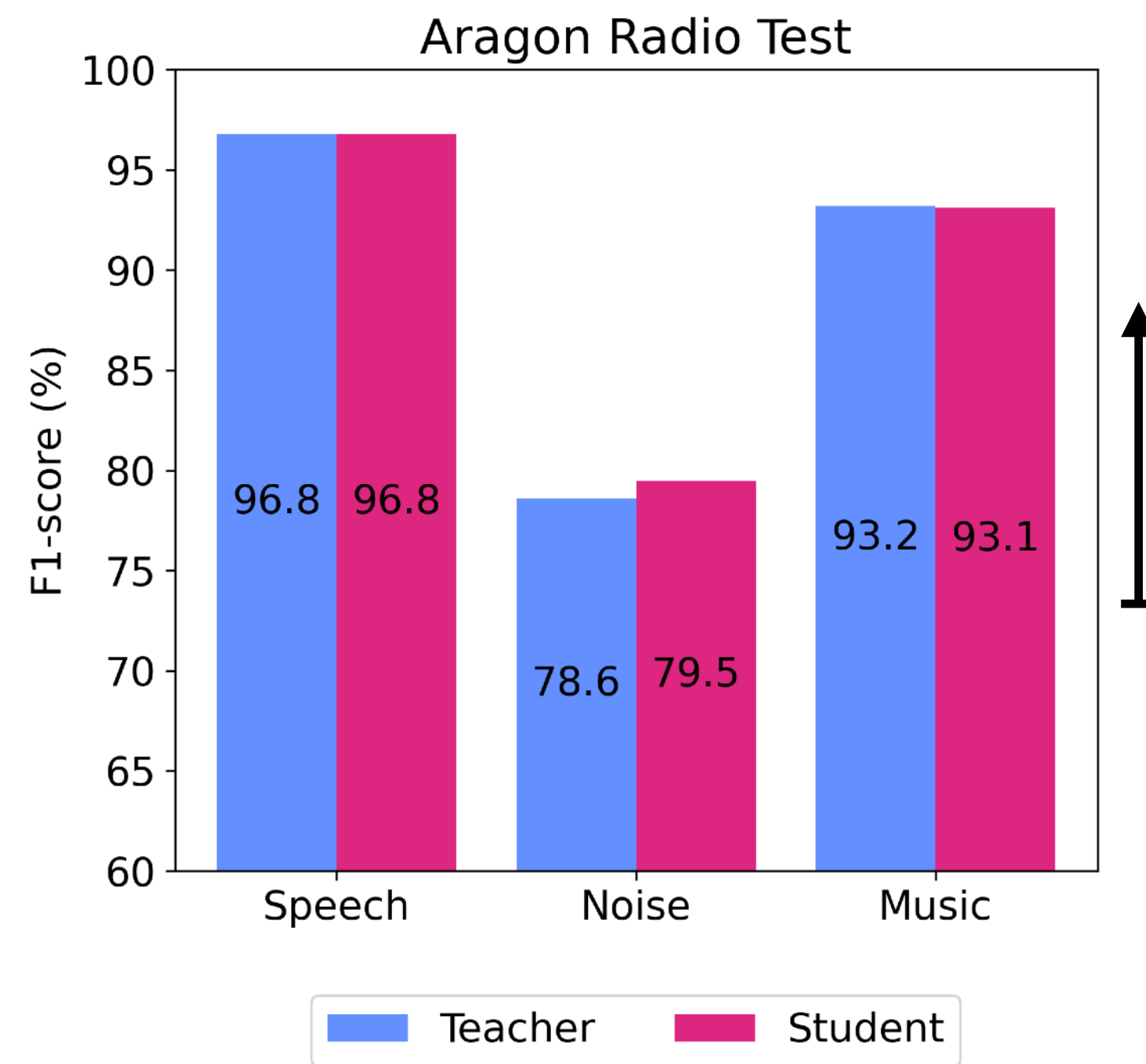
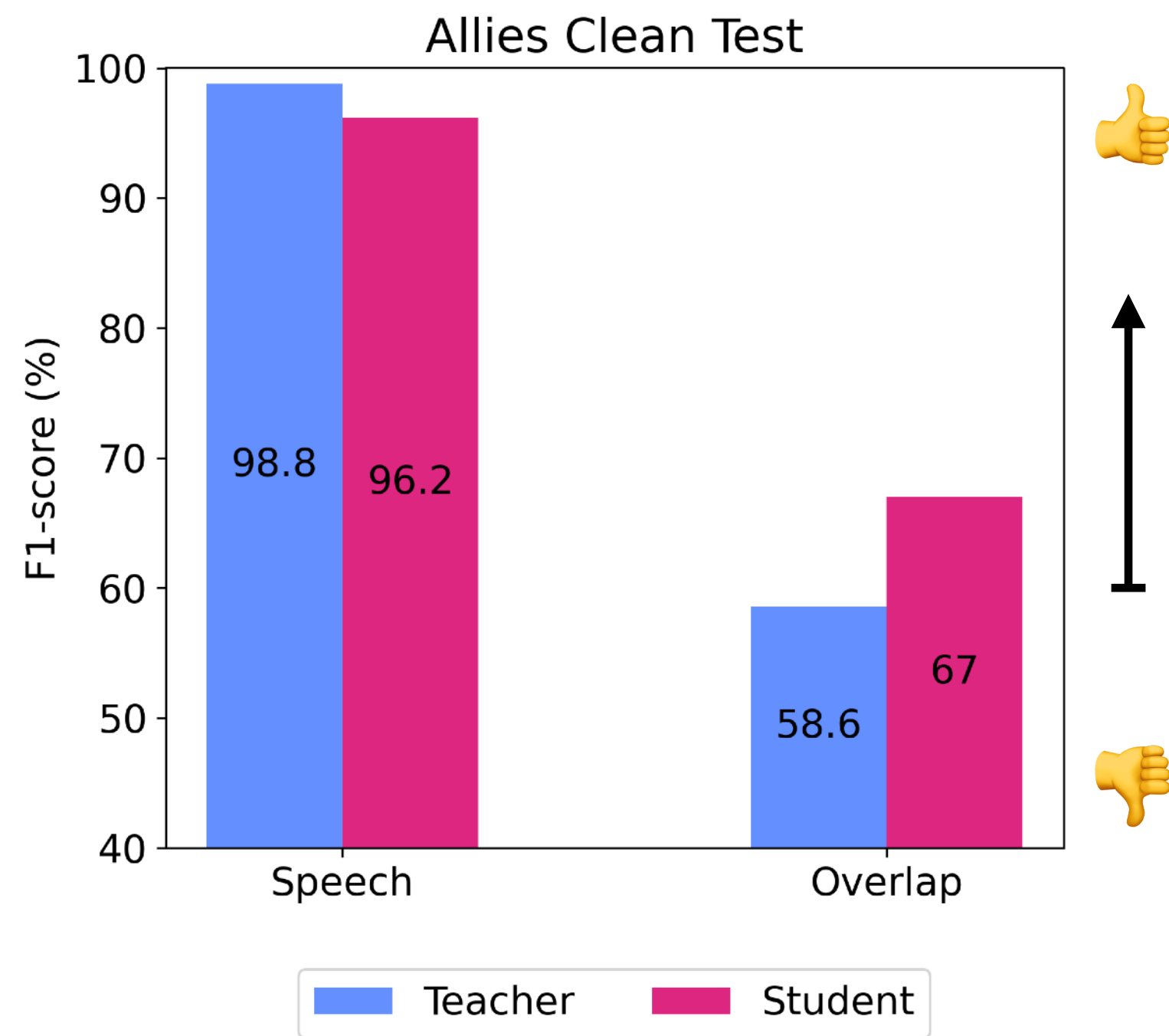
Performances de segmentation

Teacher vs. Student



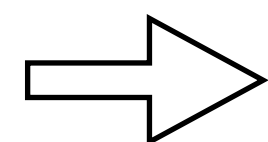
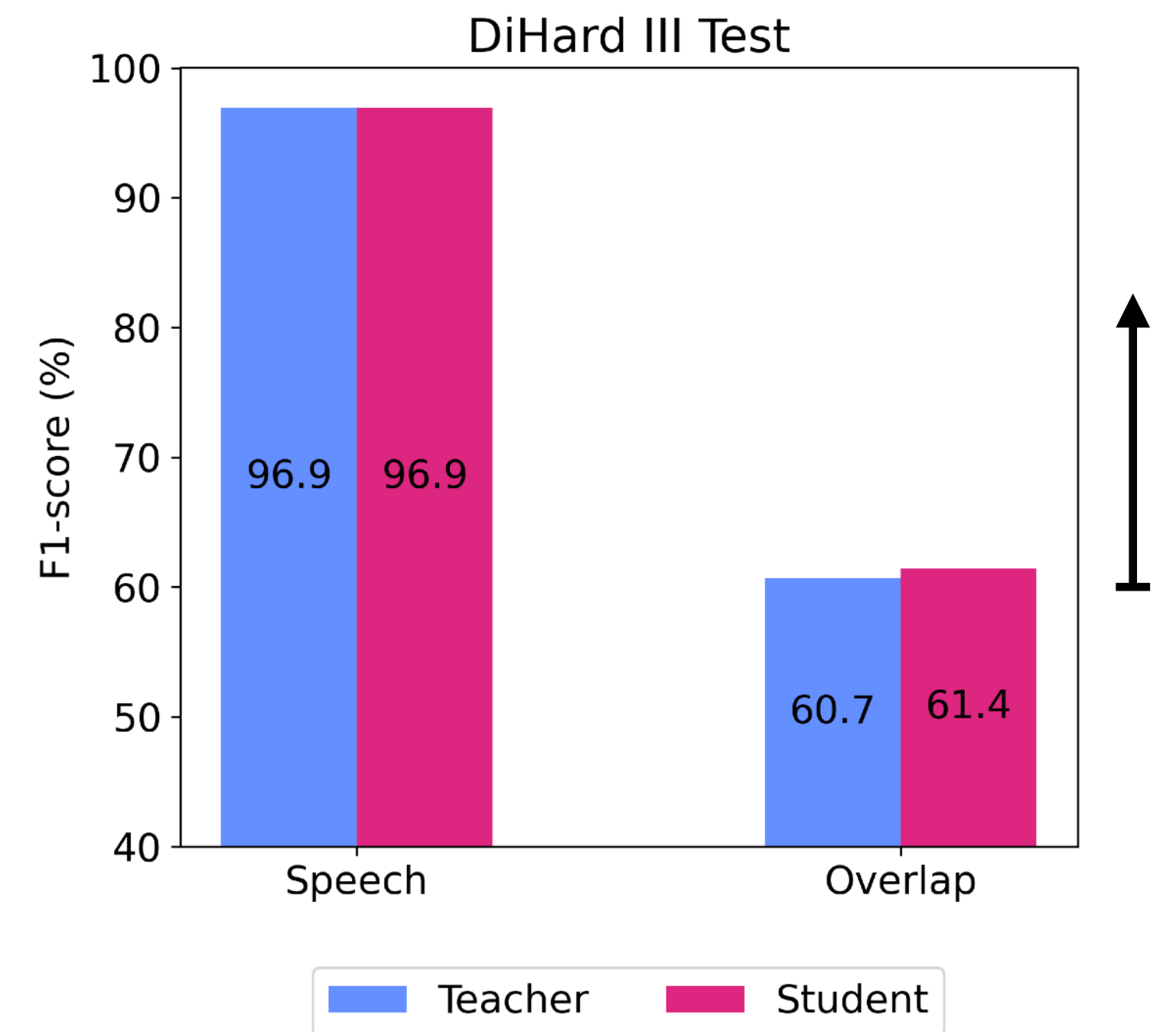
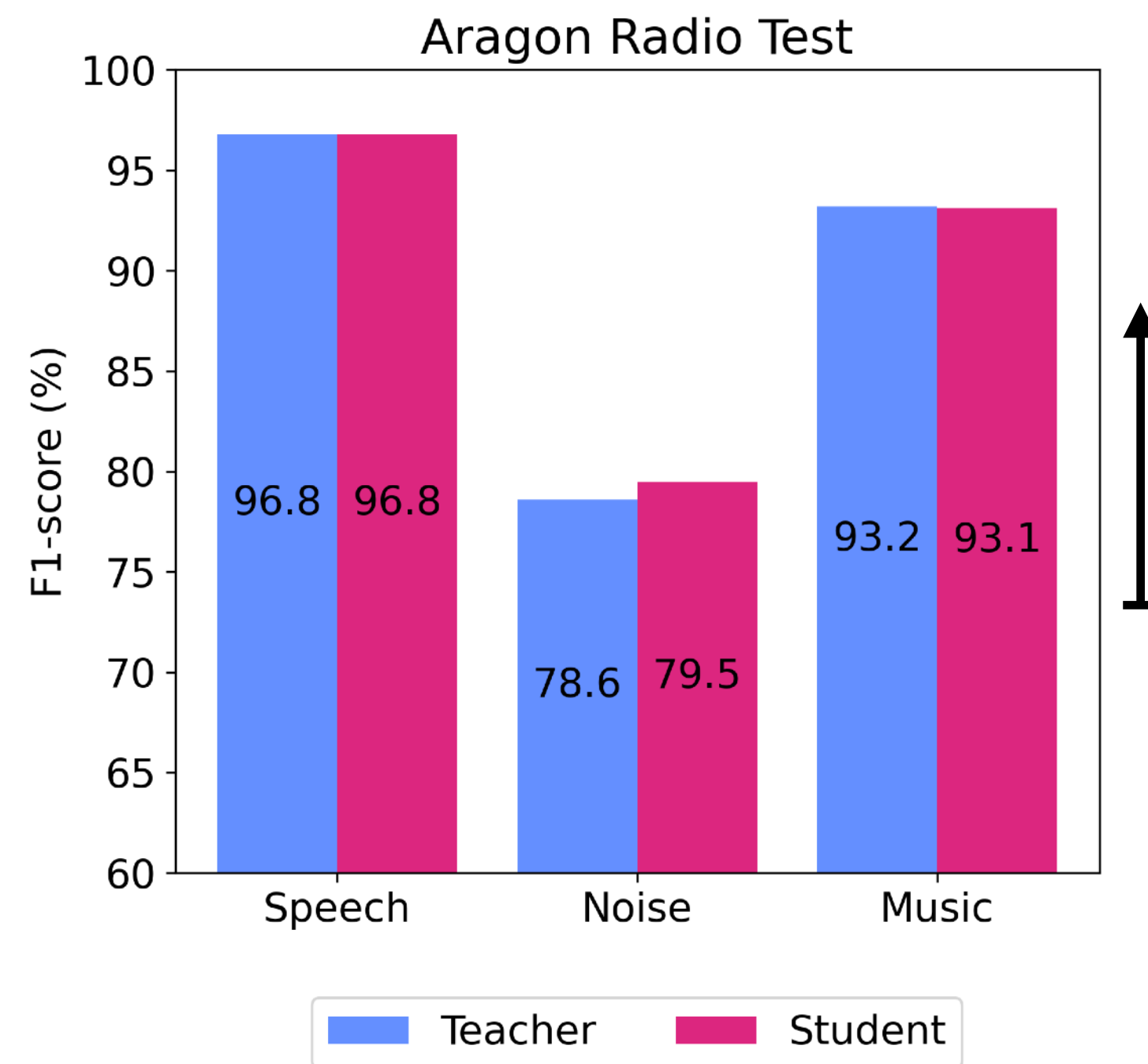
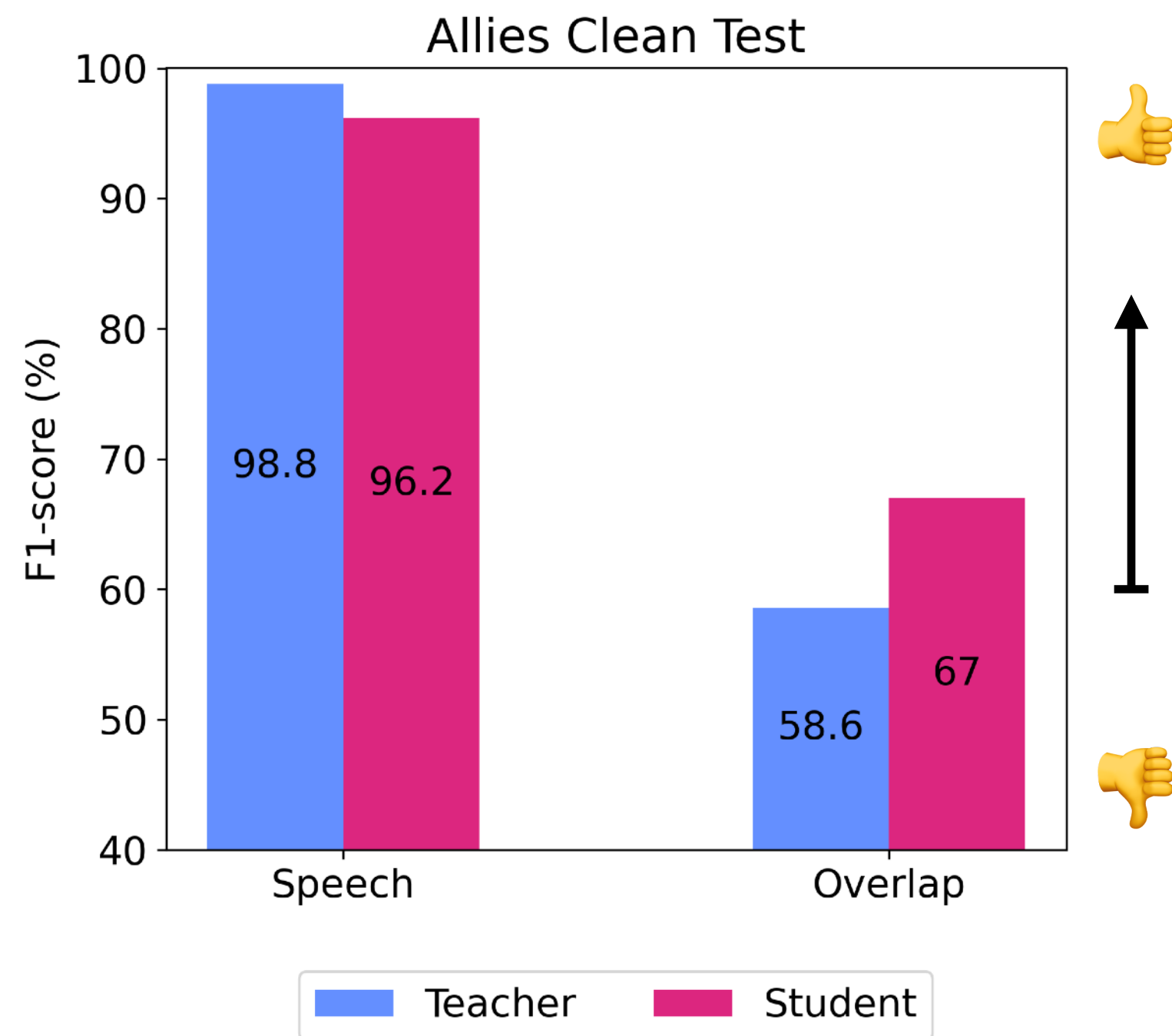
Performances de segmentation

Teacher vs. Student



Performances de segmentation

Teacher vs. Student



Le *student* offre les mêmes performances que le modèle original

Sommaire

1. Conception d'un système de segmentation audio explicable

1.1. Système *teacher*

1.2. Factorisation matricielle non-négative

1.3. Système *student* explicable

2. Évaluation des performances

2.1. Protocole

2.2. Résultats

3. Extraction d'explications

3.1. Objectif d'explication

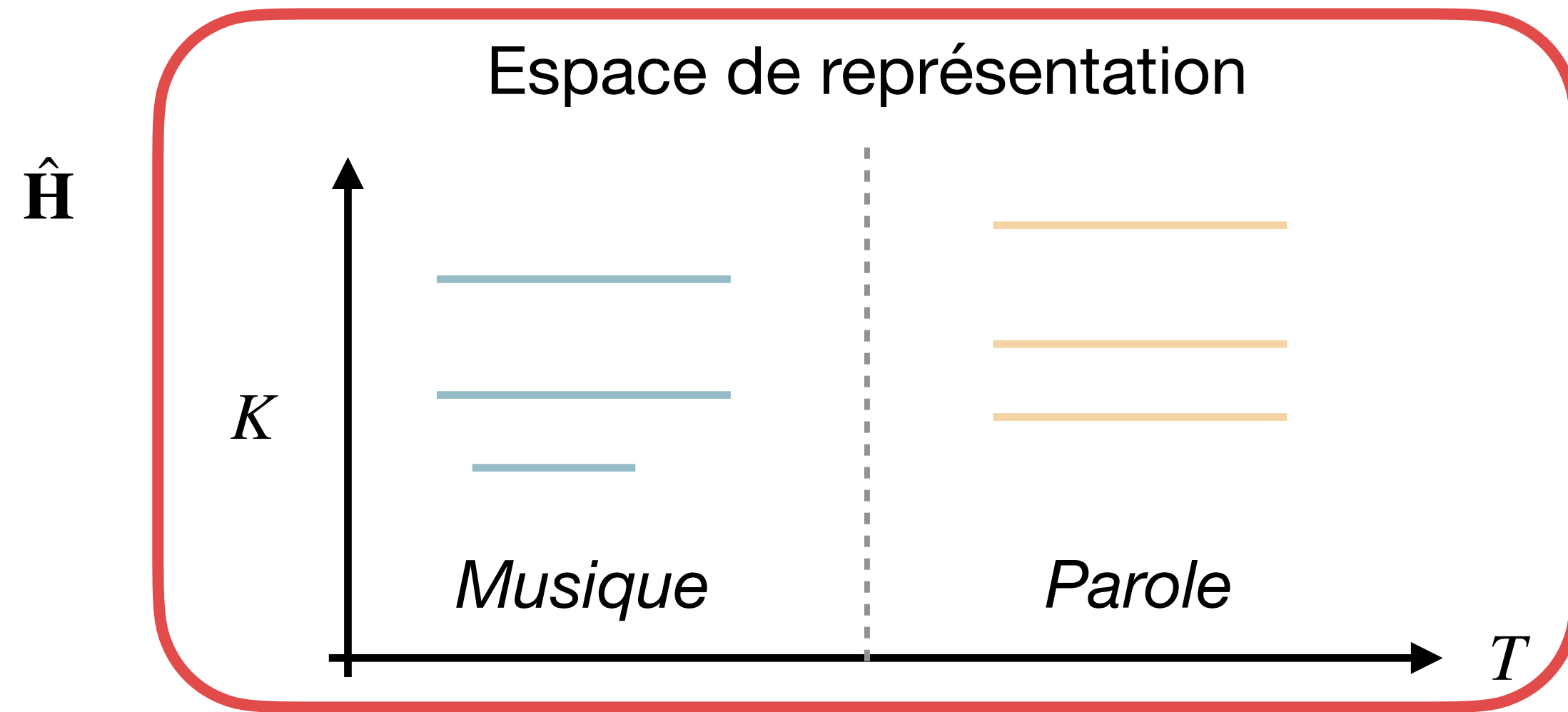
3.2. Sélection des composantes pertinentes

3.3. Analyse des explications

4. Conclusions et perspectives

Extraction d'explications

Passage de l'espace de représentation aux fréquences

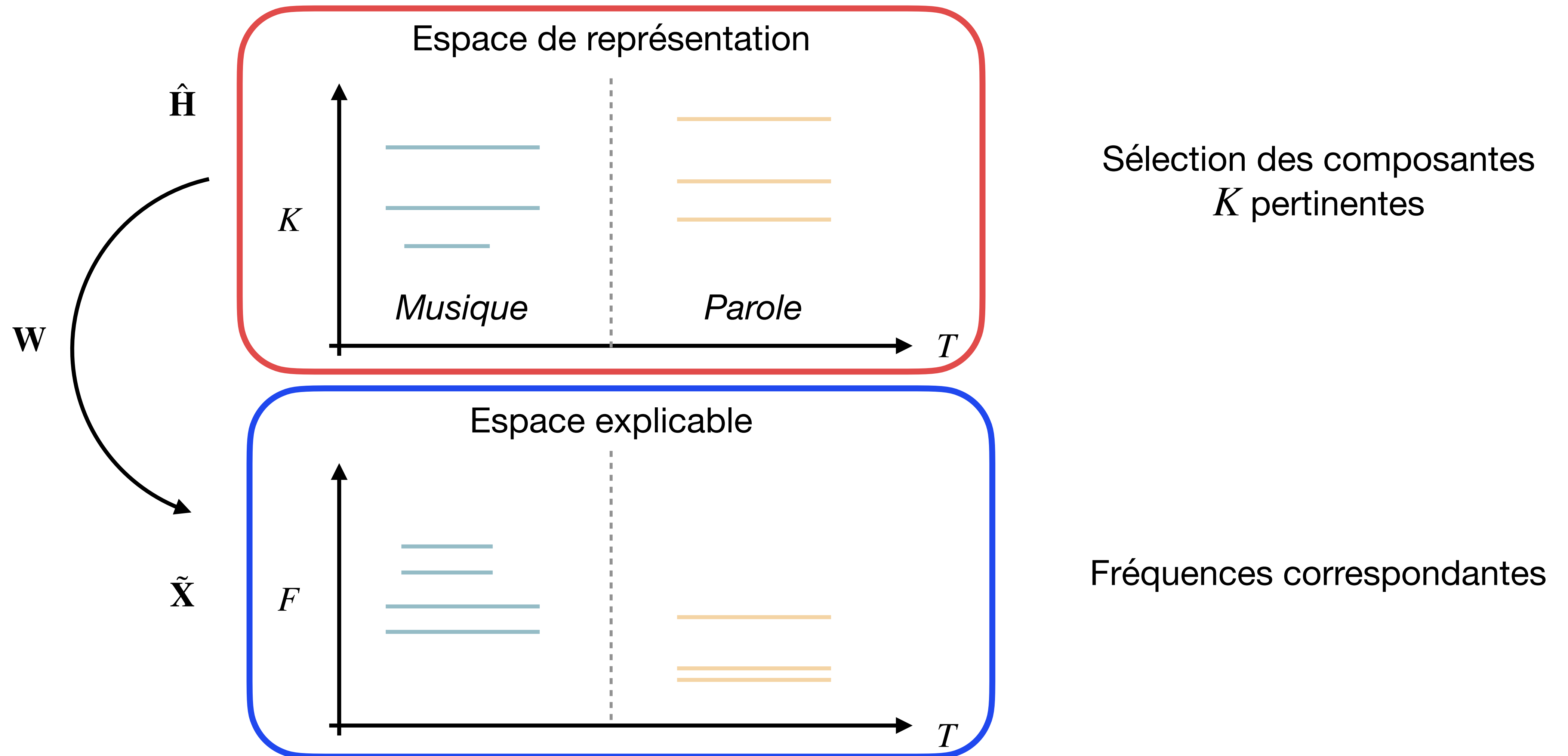


Sélection des composantes
 K pertinentes

\tilde{X}

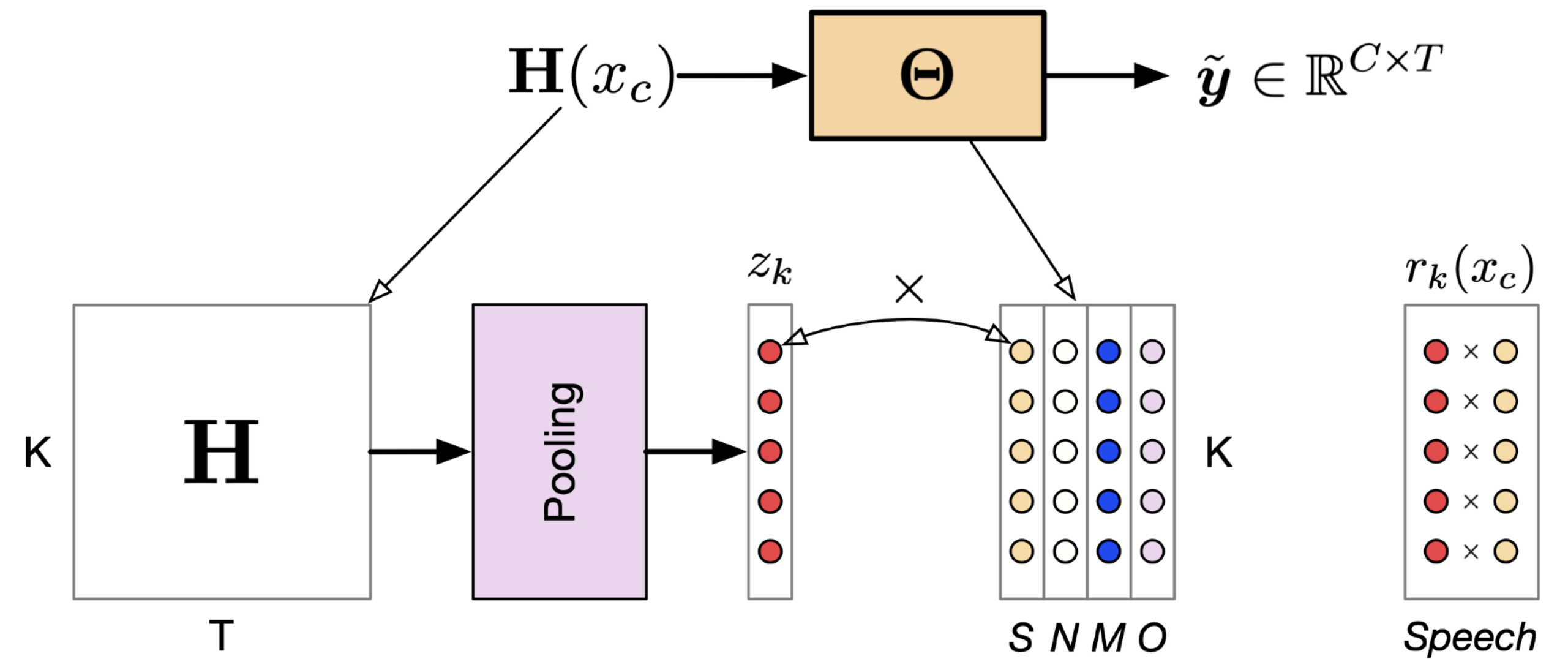
Extraction d'explications

Passage de l'espace de représentation aux fréquences



Extraction d'explications

Protocole

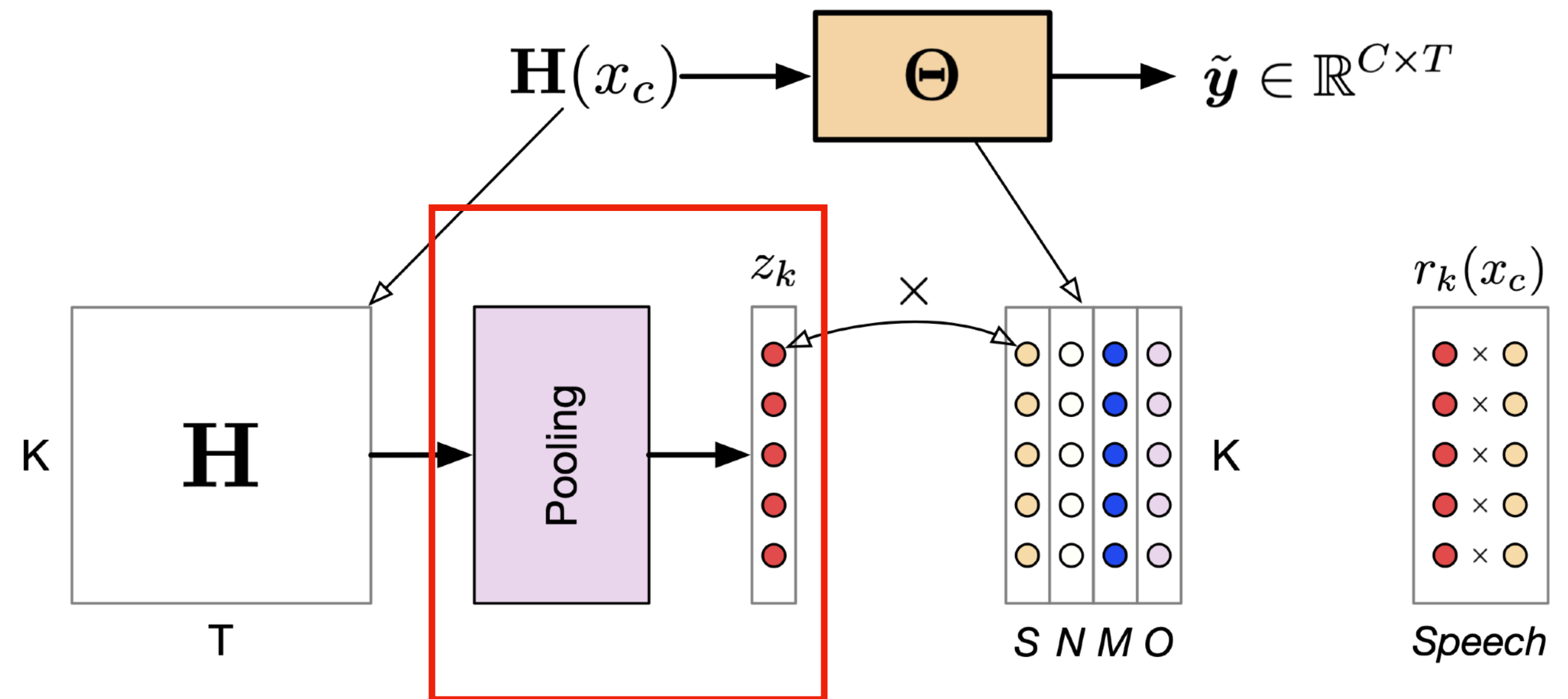


Extraction d'explications

Protocole

Moyenne temporelle de \mathbf{H} :

$$z_k(x_c) = \frac{1}{T} \sum_t H_{k,t}(x_c)$$

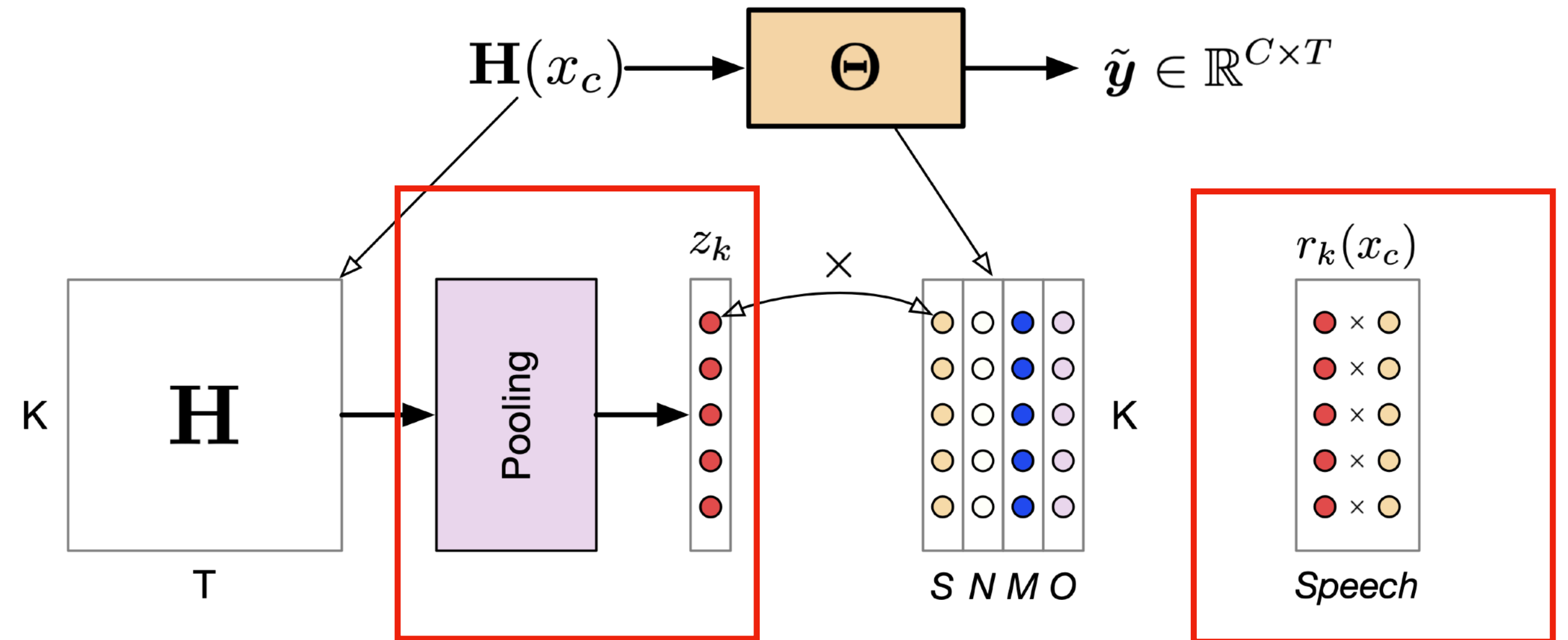


Extraction d'explications

Protocole

Moyenne temporelle de \mathbf{H} :

$$z_k(x_c) = \frac{1}{T} \sum_t H_{k,t}(x_c)$$



Mesure de la pertinence des composantes NMF pour détecter une classe c :

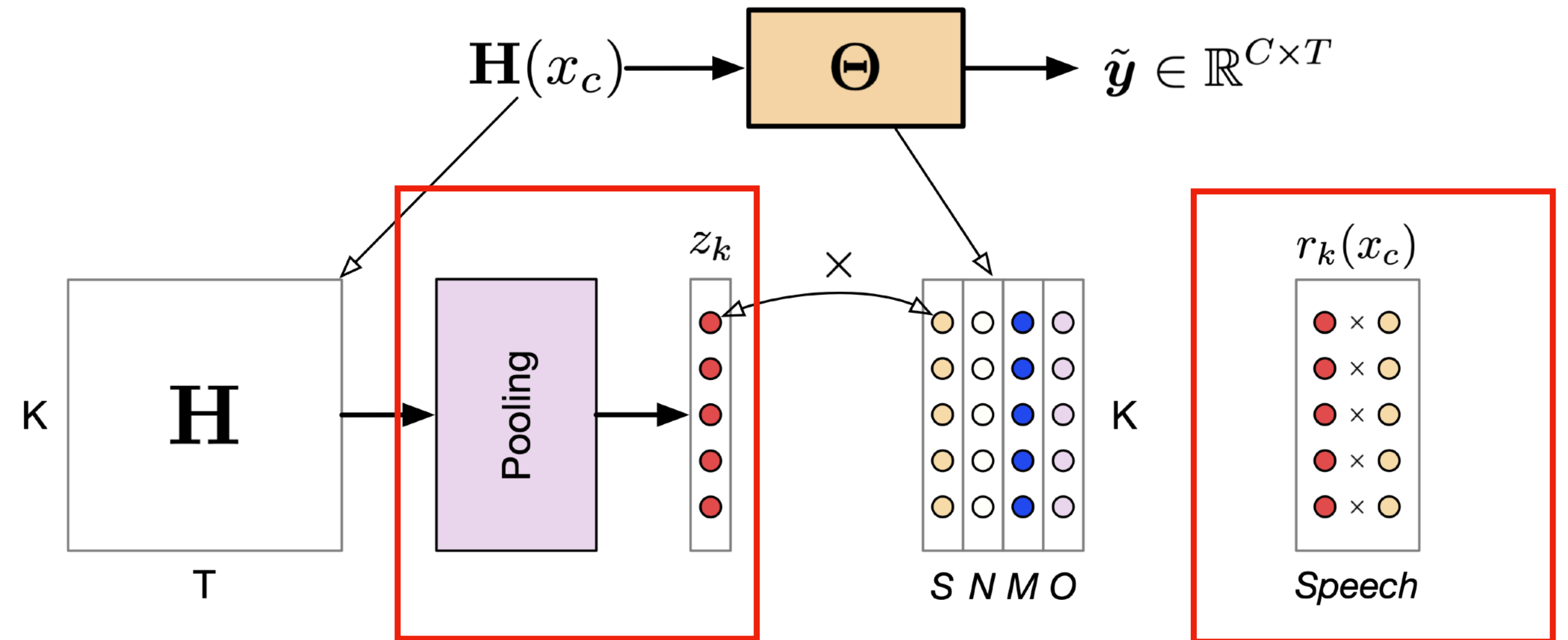
$$r_k(x_c) = z_k(x_c) \times \theta_{k,c}$$

Extraction d'explications

Protocole

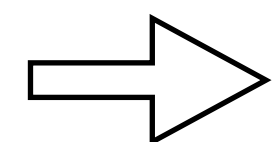
Moyenne temporelle de \mathbf{H} :

$$z_k(x_c) = \frac{1}{T} \sum_t H_{k,t}(x_c)$$



Mesure de la pertinence des composantes NMF pour détecter une classe c :

$$r_k(x_c) = z_k(x_c) \times \theta_{k,c}$$



Un seuil peut être appliqué à r_k afin de sélectionner les composantes NMF

Extraction d'explications

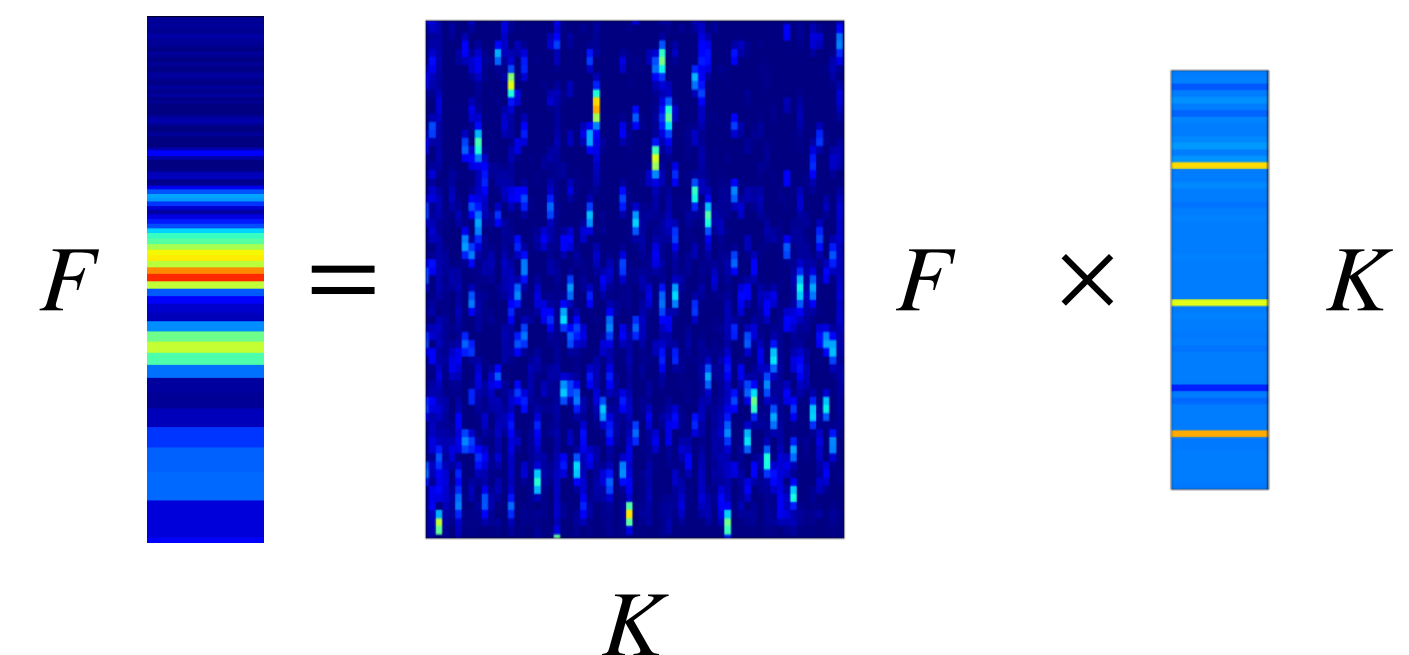
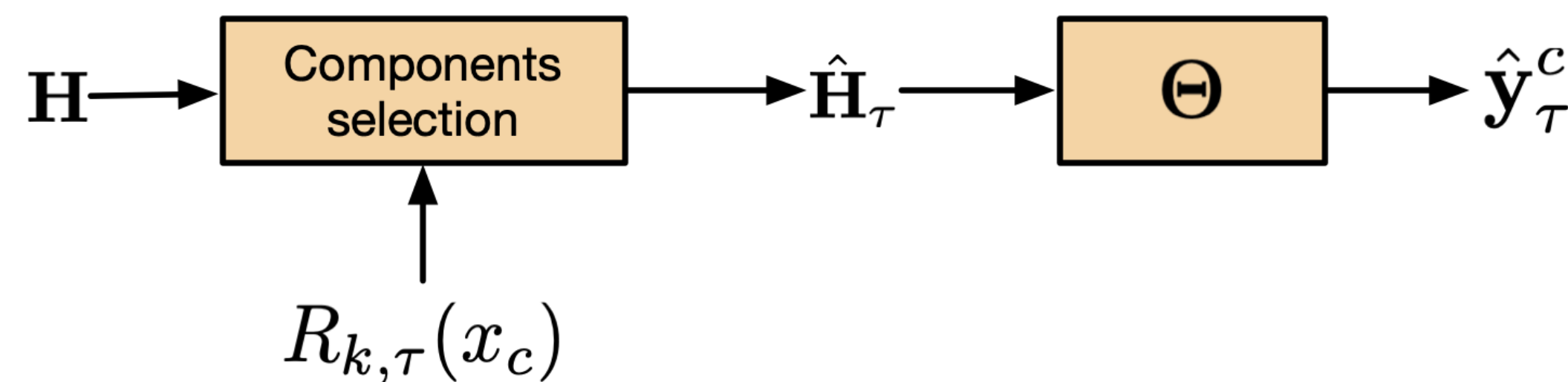
Protocole

Rappel: \mathbf{W} permet projeter les activations \mathbf{H} dans le domaine des fréquences

$$\tilde{\mathbf{X}} = \mathbf{W}\mathbf{H}(x)$$

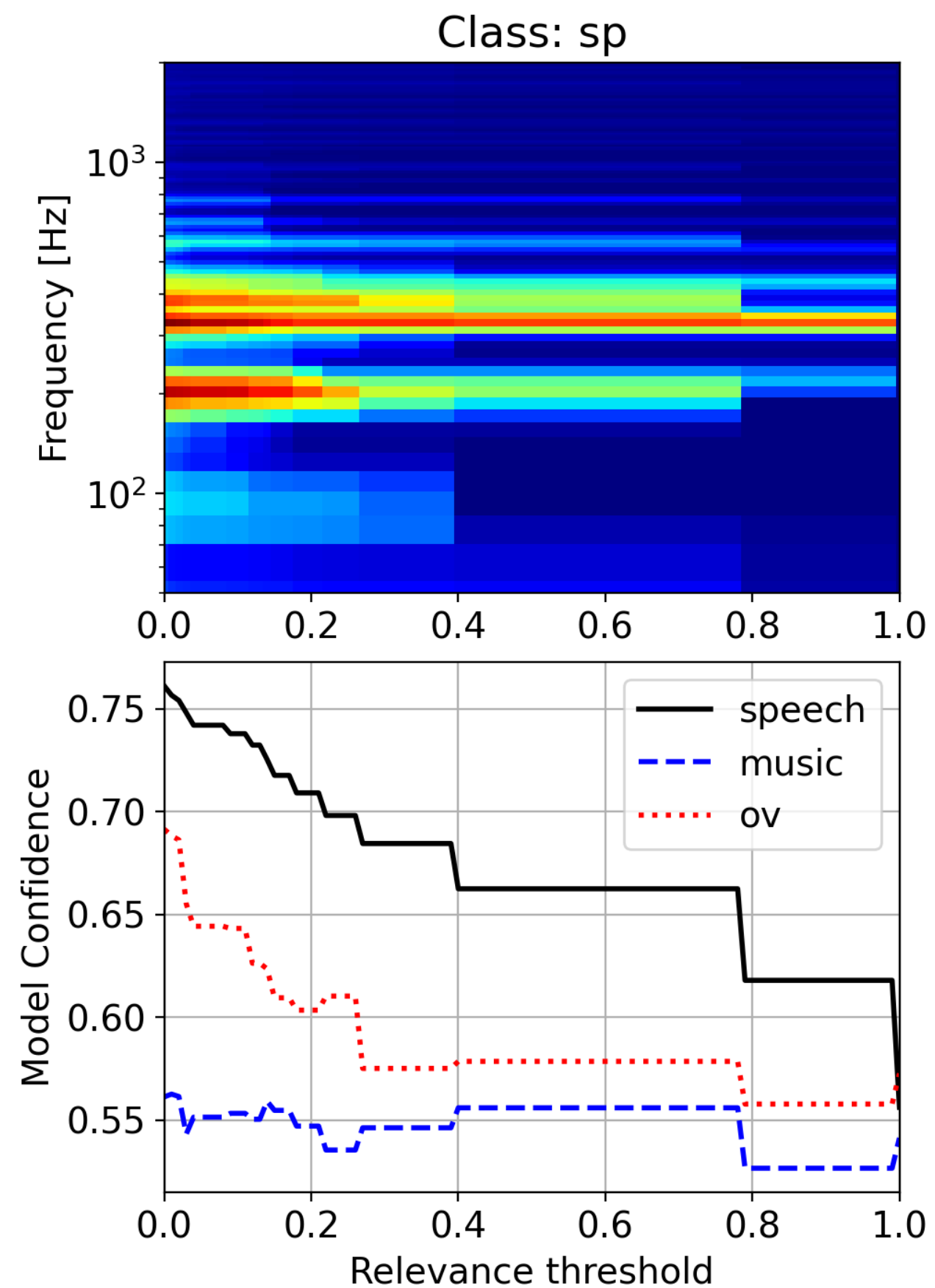
💡 Sélection des composantes pertinentes puis projection afin de le visualiser dans un espace connu

$$R_{k,c}(\tau) = \begin{cases} r_{k,c} & \text{if } r_{k,c} > \tau \\ 0 & \text{otherwise.} \end{cases} \quad \Rightarrow \quad \tilde{\mathbf{X}}_c = \mathbf{W}\mathbf{R}_\tau(x_c)$$



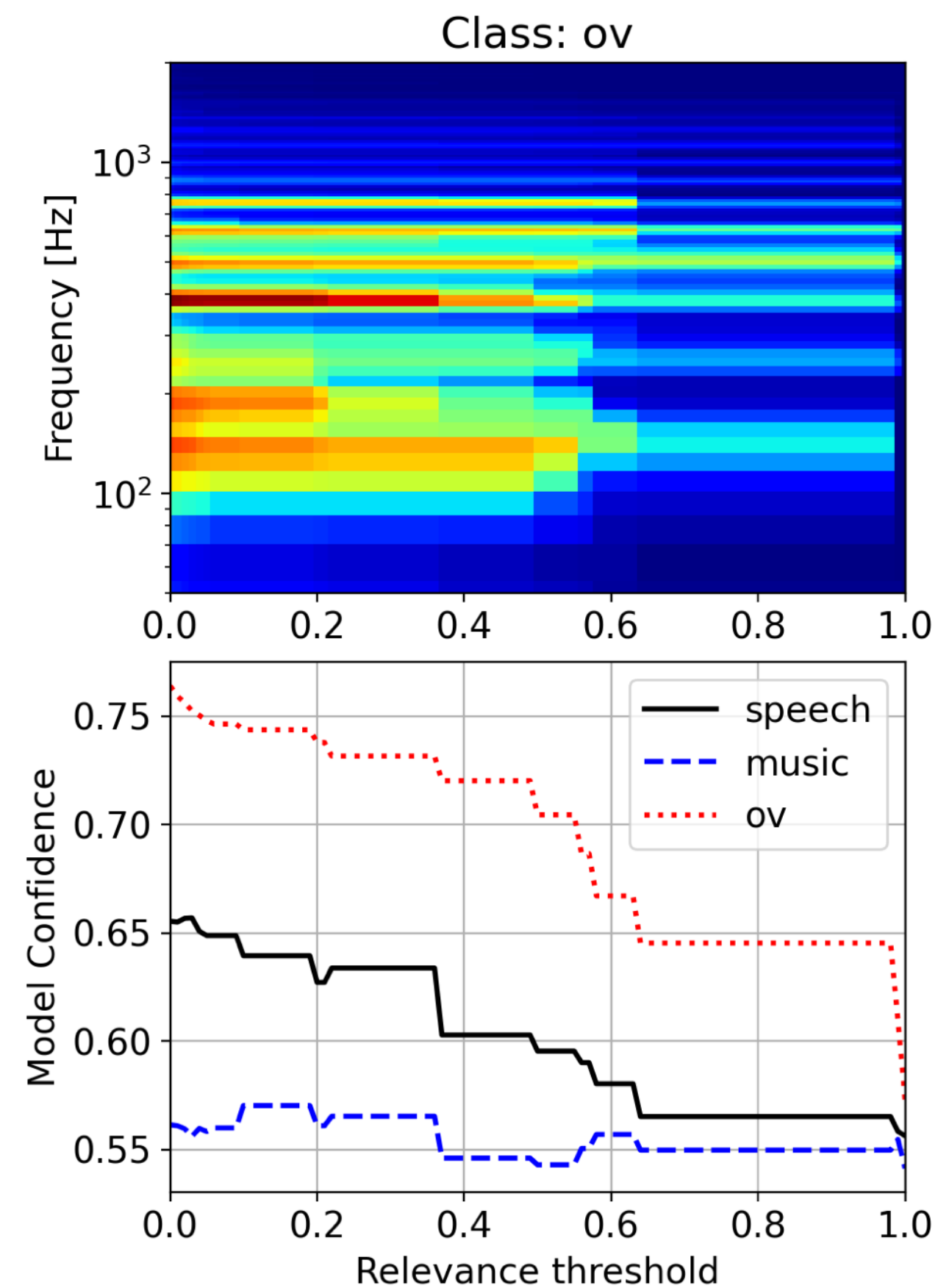
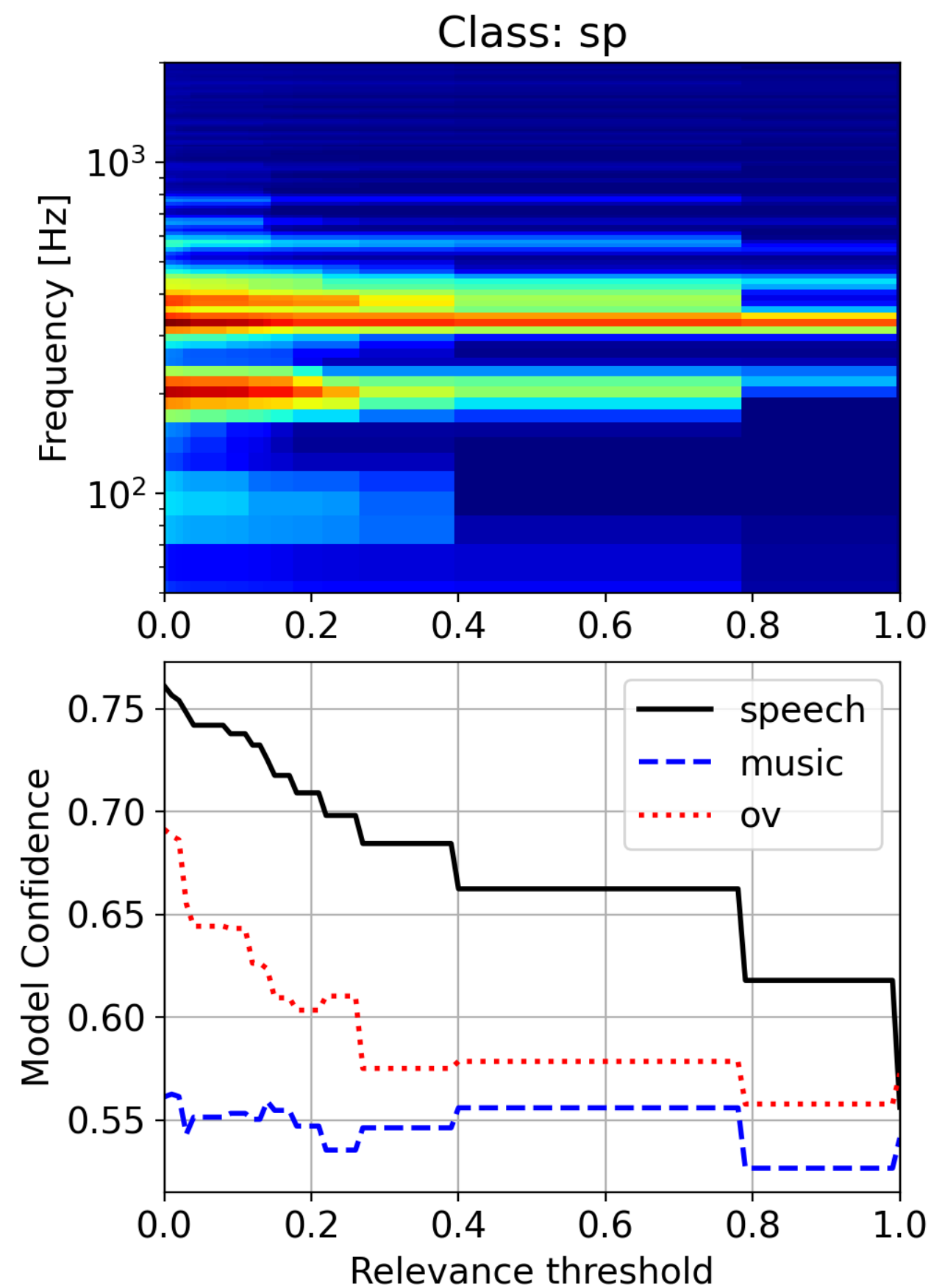
Extraction d'explications

Explications locales : quelles fréquences sont utiles à la détection ?



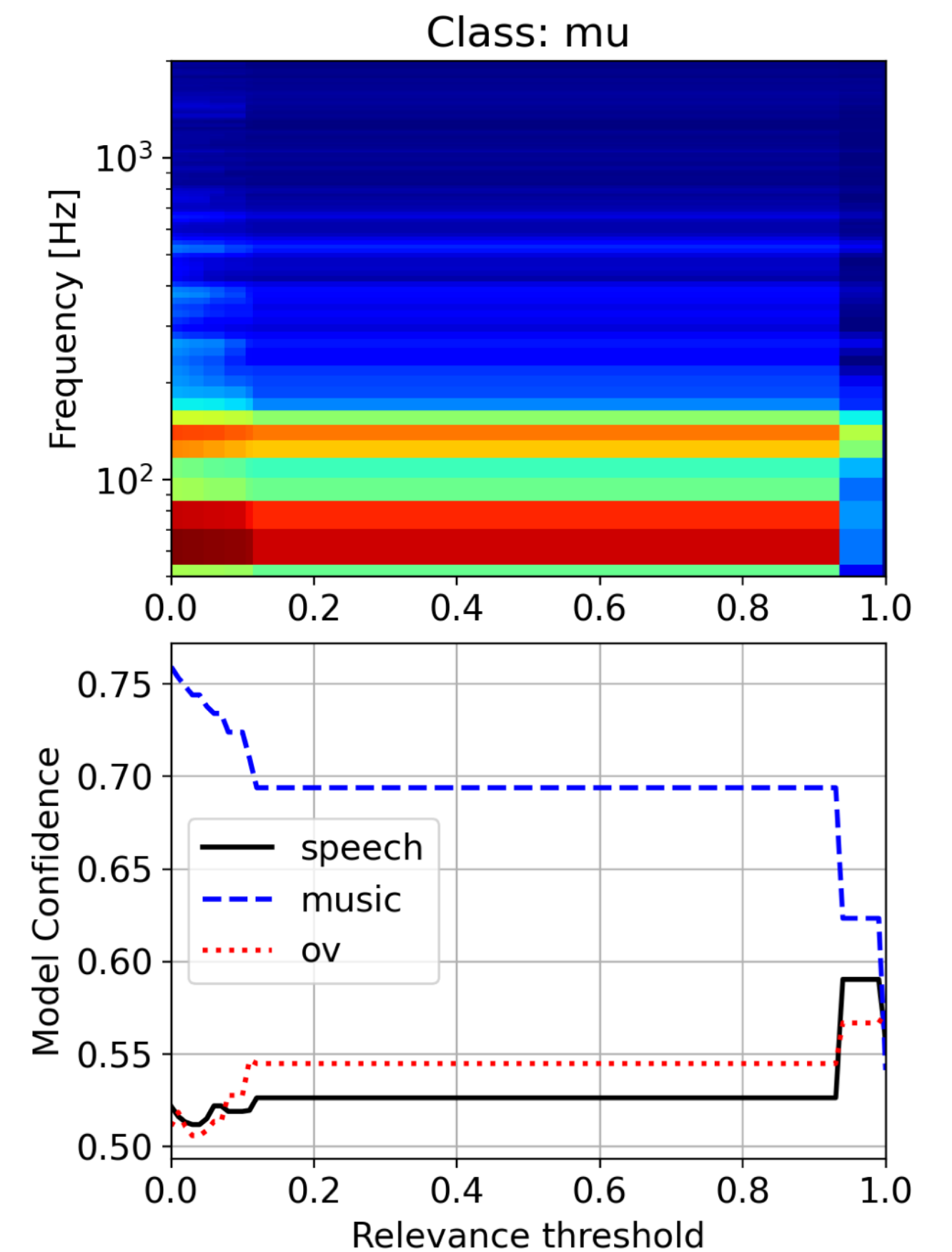
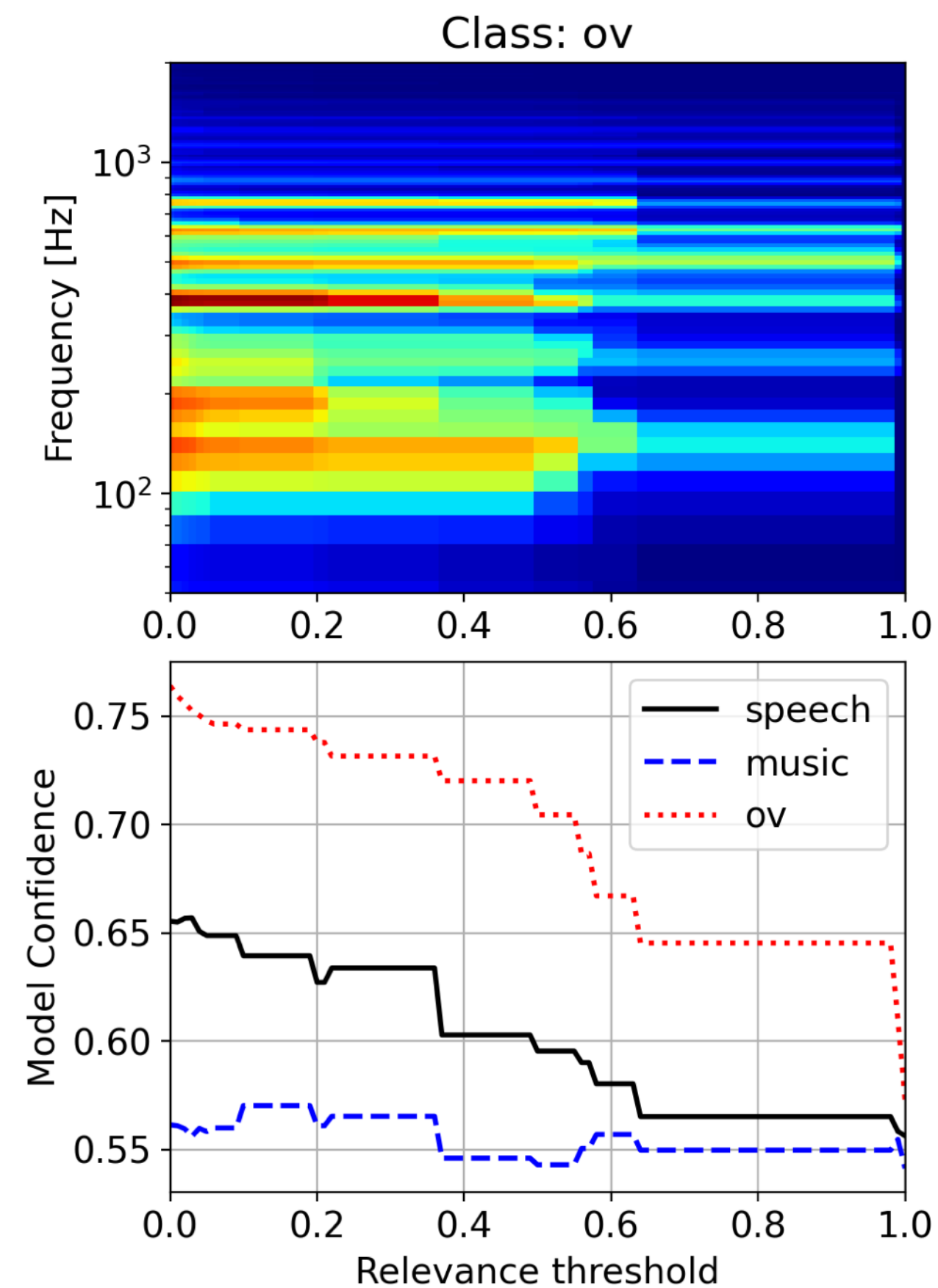
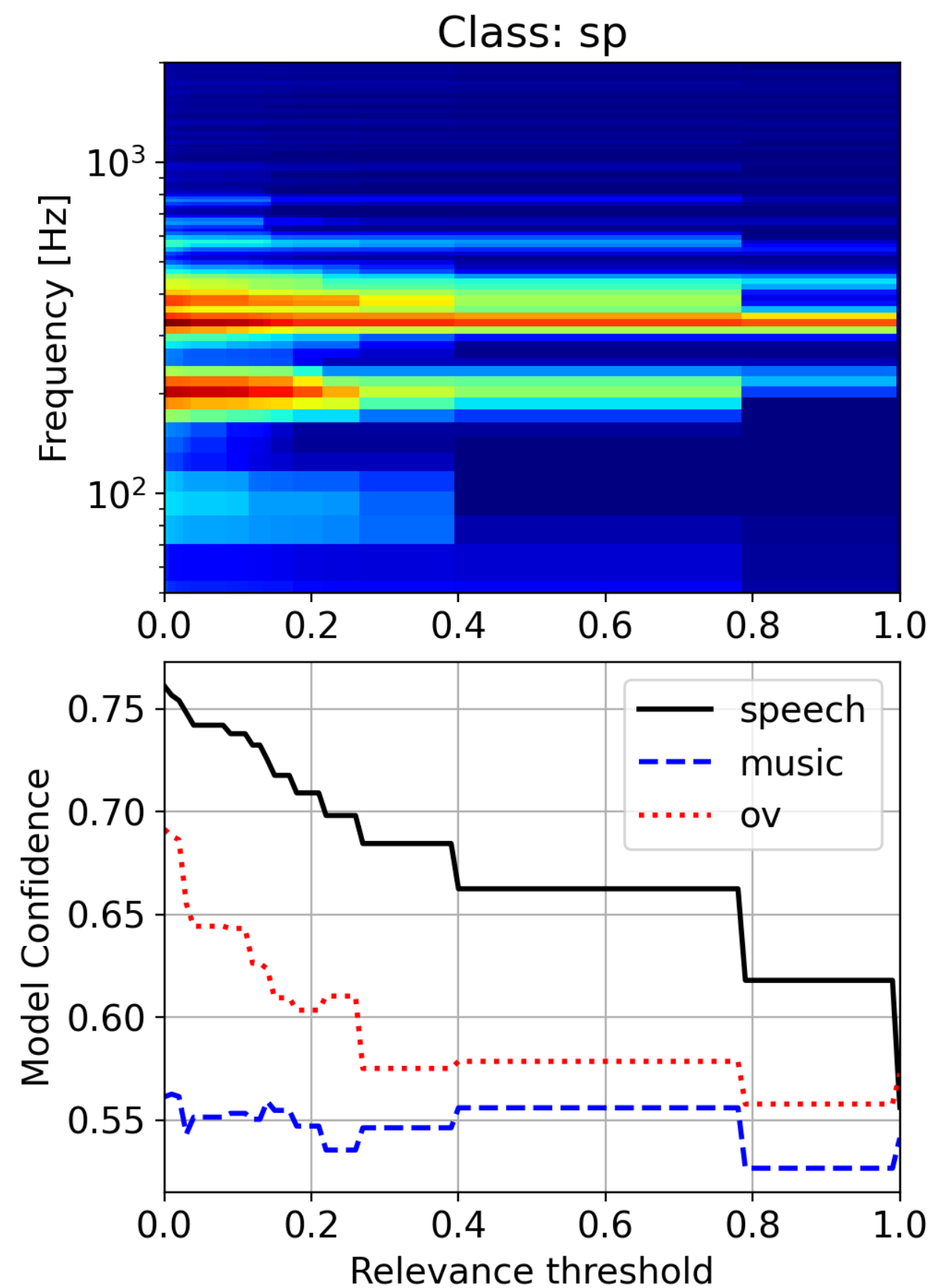
Extraction d'explications

Explications locales : quelles fréquences sont utiles à la détection ?



Extraction d'explications

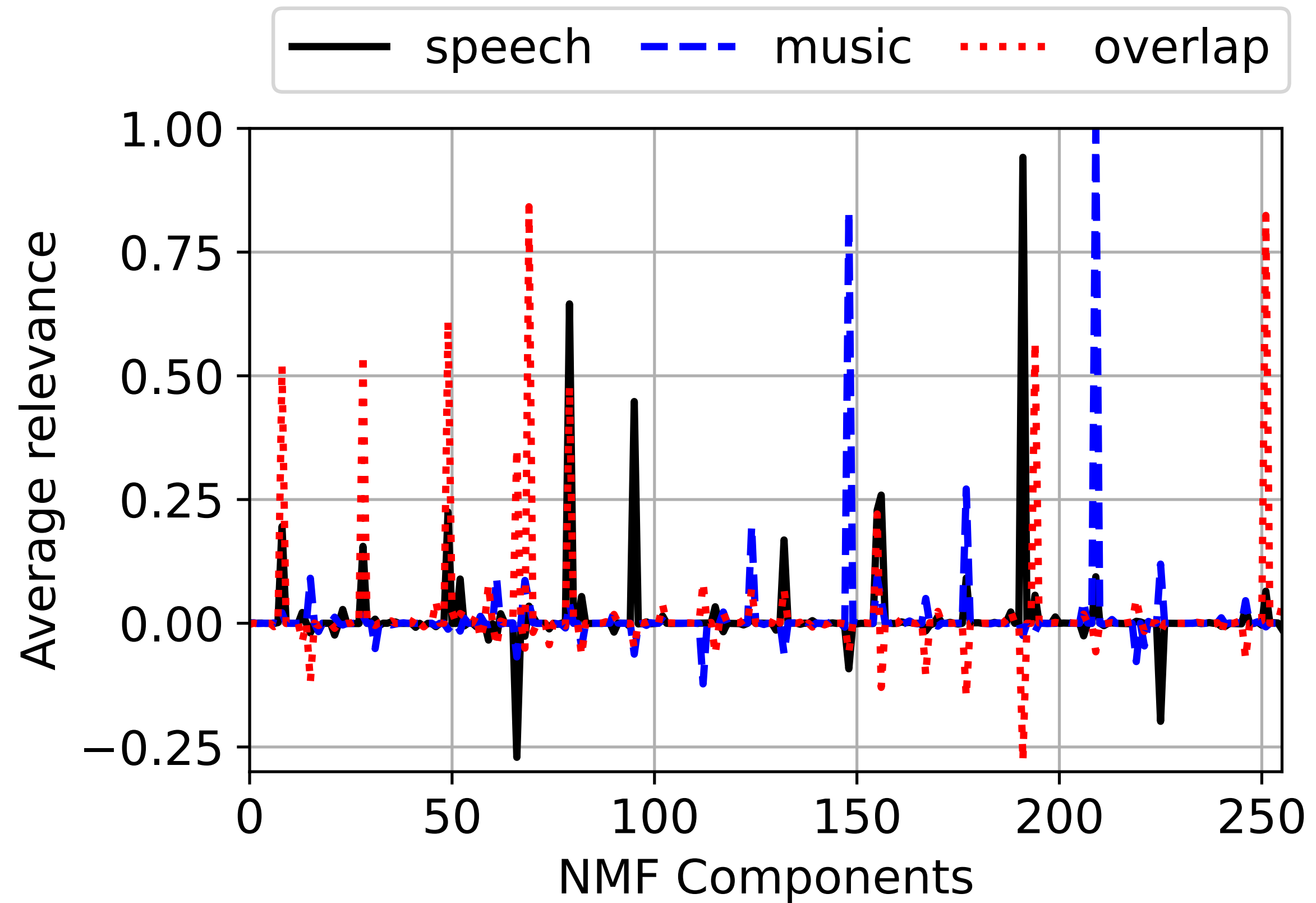
Explications locales : quelles fréquences sont utiles à la détection ?



Extraction d'explications

Explications globales : quel est l' « ADN » des classes ?

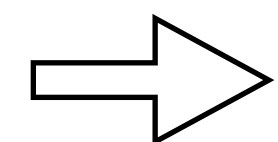
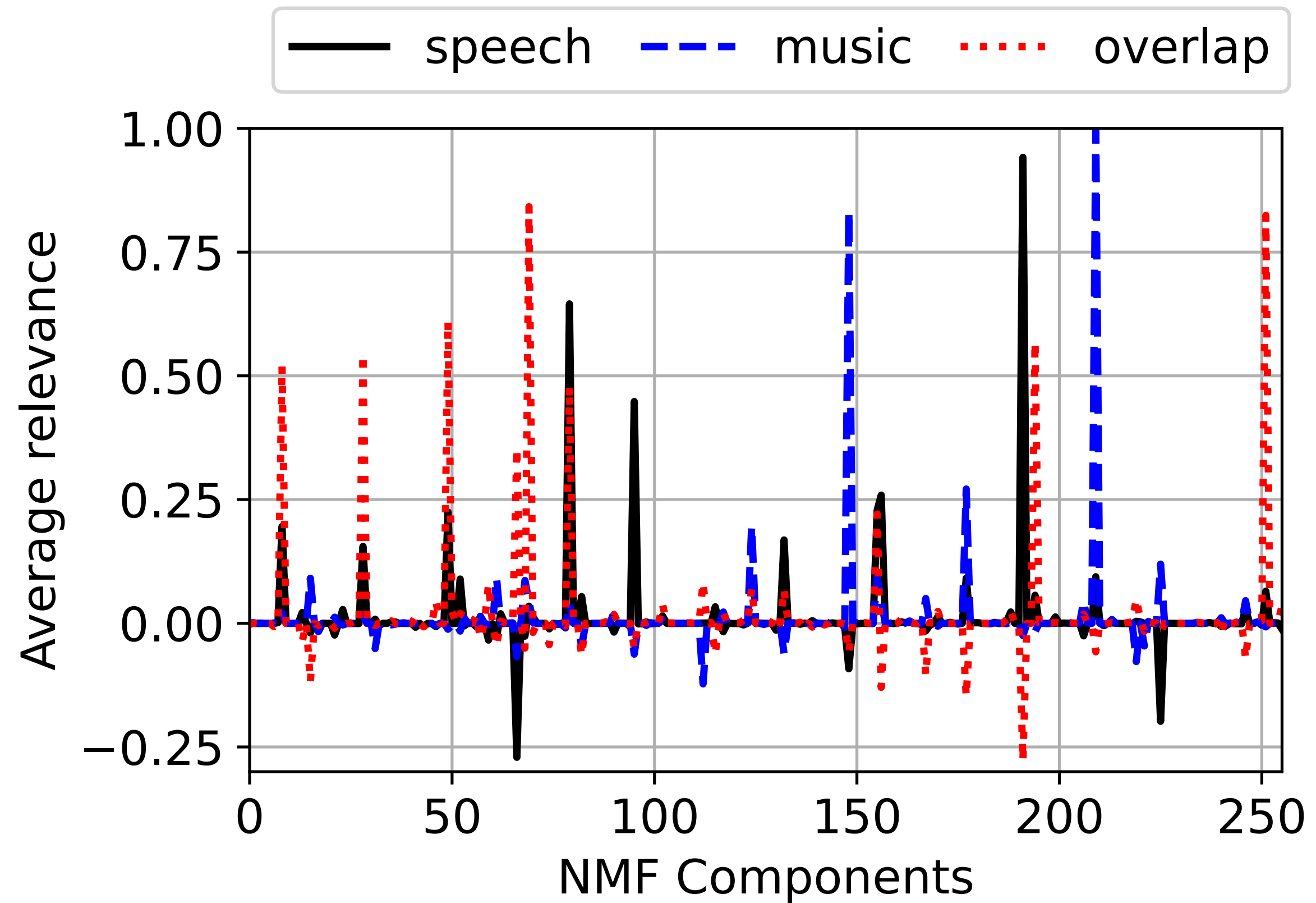
$$\mathbf{p}_c = \frac{1}{N_{\mathcal{D}}} \sum_{x_c \in \mathcal{D}} \mathbf{r}_c(x)$$



Extraction d'explications

Explications globales : quel est l' « ADN » des classes ?

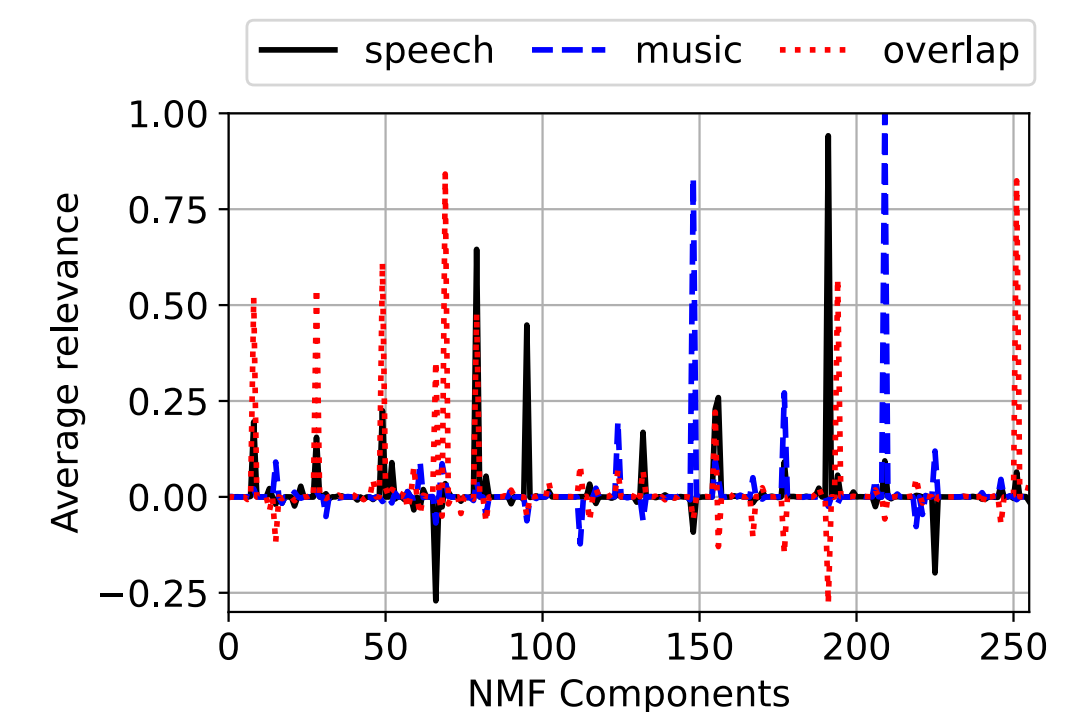
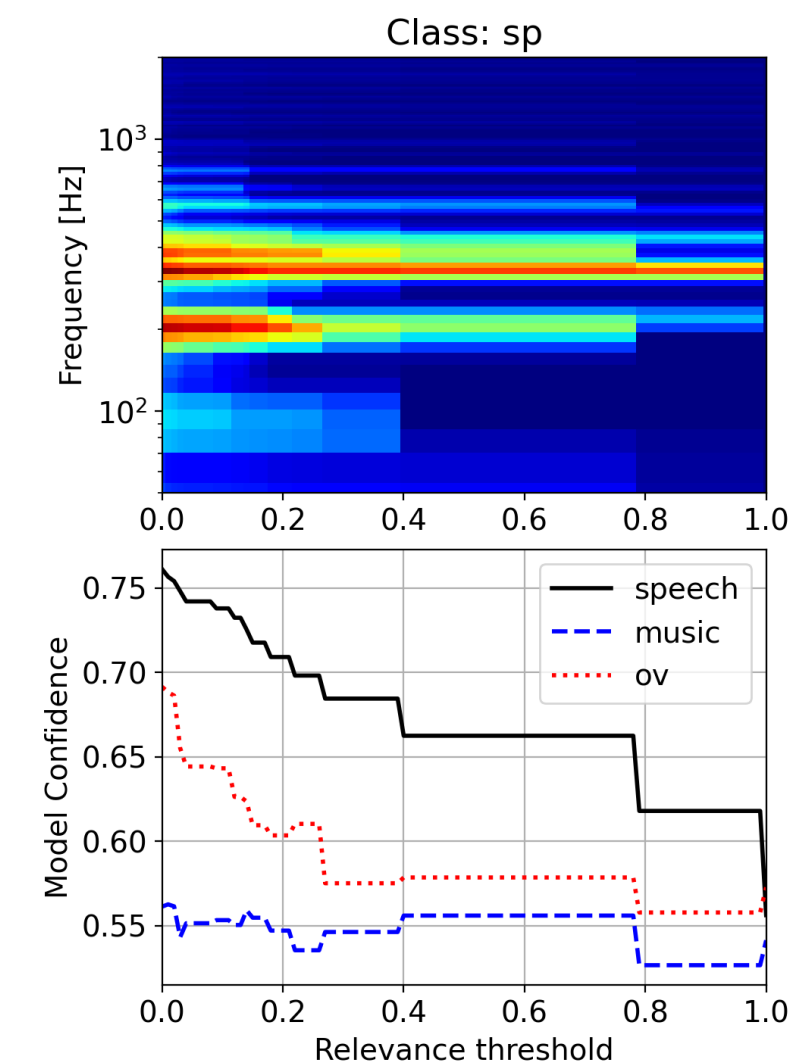
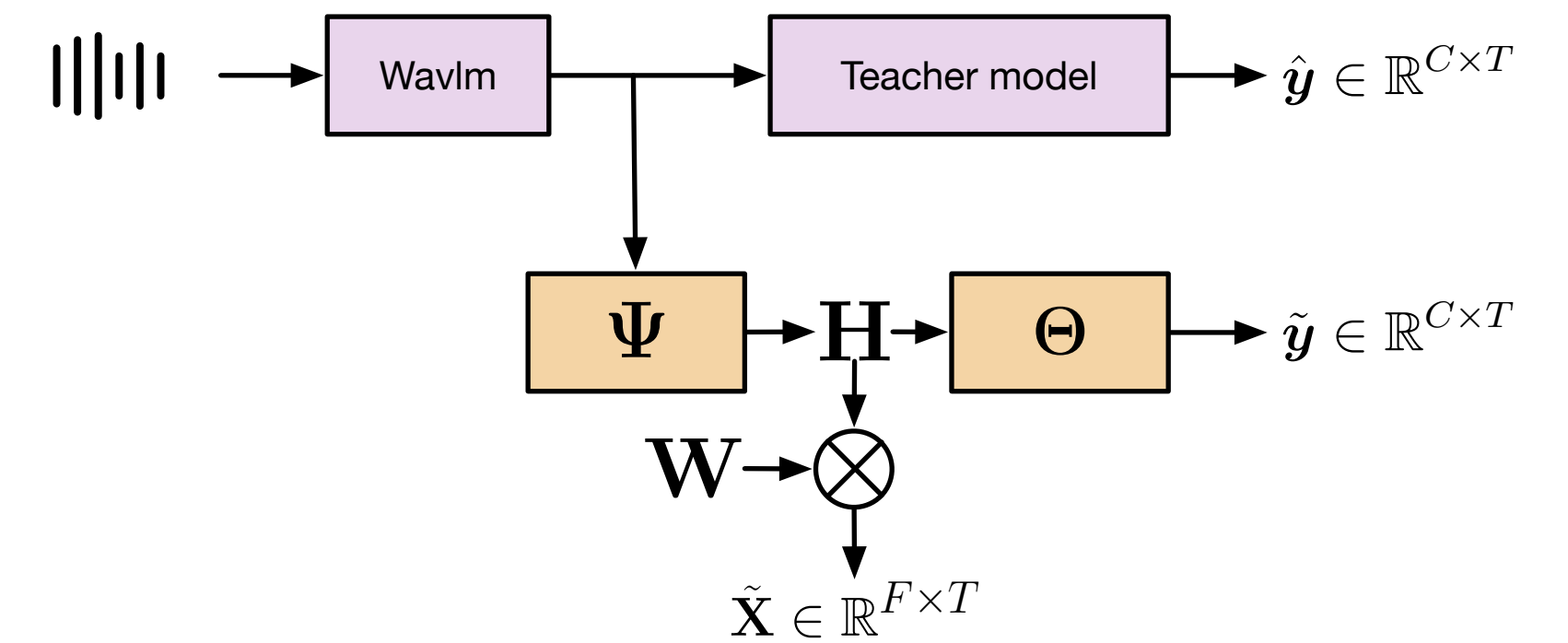
$$\mathbf{p}_c = \frac{1}{N_{\mathcal{D}}} \sum_{x_c \in \mathcal{D}} \mathbf{r}_c(x)$$



La pertinence moyenne pour chaque classe identifie les composantes représentatives

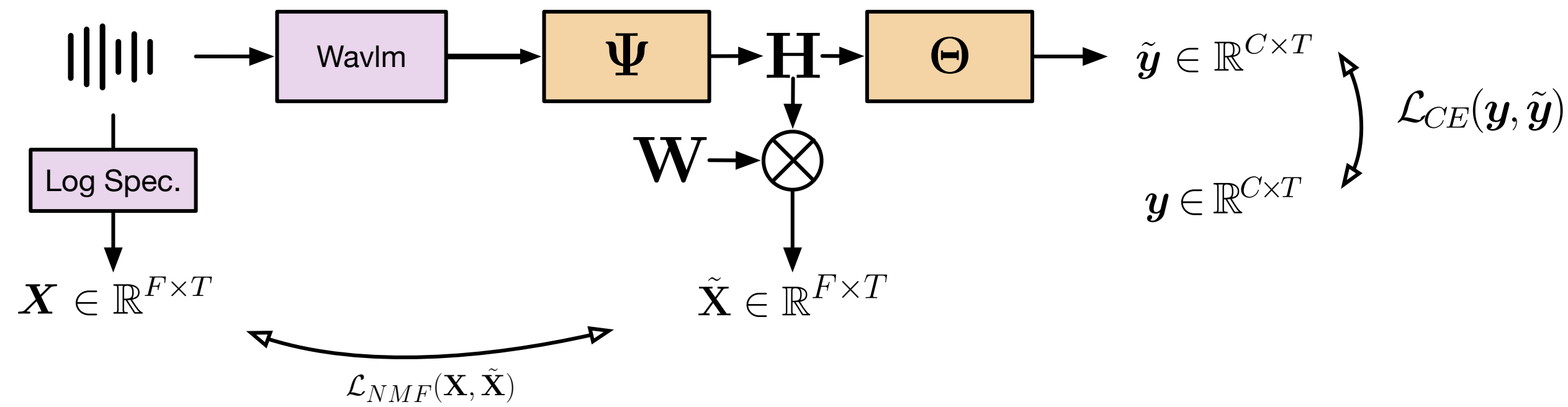
Conclusions

- Modèle de segmentation audio (parole, parole superposée, musique et bruit)
- Apprentissage d'un système explicable par distillation de connaissances
- Ce système utilise la NMF pour projeter l'espace de représentation dans le domaine des fréquences
- Mêmes performances que le *teacher*
- Possibilité d'extraire des explications locales (pour un exemple) et globales (pour un jeu de données)



Perspectives

Suppression du *teacher* pour définir un modèle explicable par construction



Probing de la matrice \mathbf{H} pour une meilleure explication :
« cette matrice n'est-elle pas sur-optimisée ? »

Détections de genre de musique, de phonèmes, d'évènements sonores...

