

# person instance graphs for person recognition in multimedia data

**Hervé Bredin (CNRS / LIMSI)**

Anindya Roy - Viet-Bac Le - Claude Barras



# outline

- context
- **building** person instance graphs
- **mining** person instance graphs
- experimental **results**
- conclusion

# the REPERE challenge

- multimodal person recognition in TV shows
- three  projects  
co-funded by 

QCOMPERE

SODA

PERCOL

**LIMSI** et al.

**LIUM** et al.

**LIF** et al.

- two evaluation campaigns in 2013 and 2014

- data collection & annotation by 
- evaluation by 

# multimodal person recognition

- **multiple** sources of information
  - **audio stream**  
speaker diarization & identification  
speech transcription (ASR)
  - **visual stream**  
face clustering & recognition  
optical character recognition (OCR)
  - **text stream** (from ASR and OCR)  
named entity detection  
name normalization



# the QCOMPERE project



# the REPERE challenge



*who speaks when?*

*who appears when?*

# the REPERE challenge



*who speaks when?*

*who appears when?*

# the holy grail

## *free annotated data*

- **dense audio annotation**
  - speaker identification
  - speech transcription
- **sparse video annotation**
  - head recognition
  - optical character recognition
- **text annotation** named entity detection (person names only)

# the holy grail

## *free annotated data*

- **dense audio annotation**

speaker identification  
speech transcription

- **sparse video annotation**

head recognition  
optical character recognition

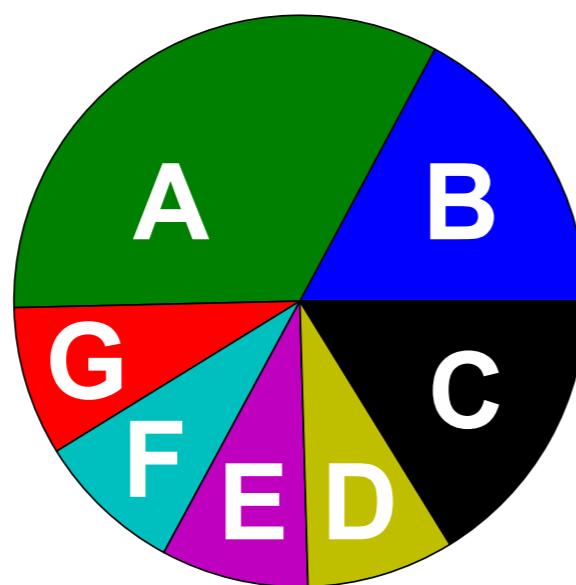
- **text annotation** named entity detection (person names only)

2013

### TRAINING SET

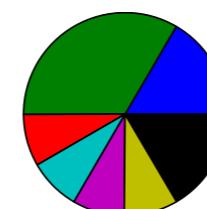
*24 hours*

- A: BFM Story
- B: LCP Info
- C: Top Questions
- D: Ça Vous Regarde
- E: Planète Showbiz
- F: Entre Les Lignes
- G: Pile Et Face



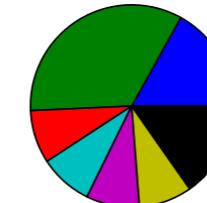
### DEVELOPMENT

*3 hours*



### TEST

*3 hours*



# the holy grail

## *free annotated data*

- **dense audio annotation**

speaker identification  
speech transcription

- **sparse video annotation**

head recognition  
optical character recognition

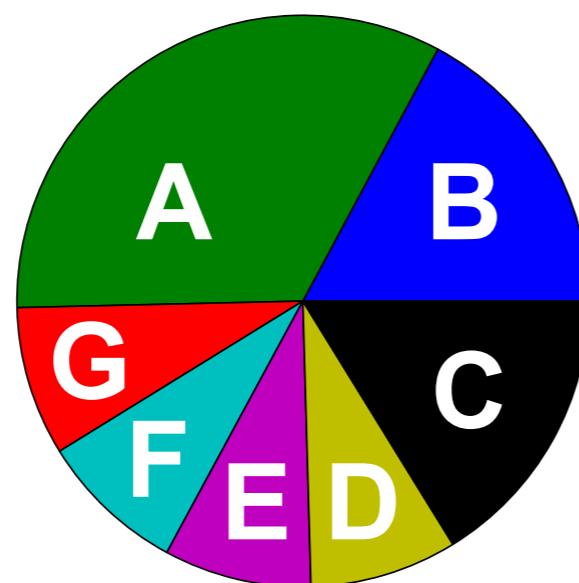
- **text annotation** named entity detection (person names only)

2013

### TRAINING SET

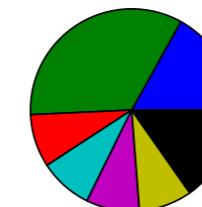
*24 hours*

- A: BFM Story
- B: LCP Info
- C: Top Questions
- D: Ça Vous Regarde
- E: Planète Showbiz
- F: Entre Les Lignes
- G: Pile Et Face



### DEVELOPMENT

*3 hours*



### TEST

*3 hours*

2014 X 2

# various subtasks

- **supervised** recognition
  - can rely** on prior biometric models (face or voice) from manually annotated data
  - can rely** on written and/or pronounced person names
- **unsupervised** recognition
  - must not rely** on prior biometric models
  - must rely** on written and/or pronounced person names
- **mono-modal** recognition
  - use audio signal only for speaker recognition
  - use facial features for face recognition
- **cross-modal** recognition
  - use written/spoken name for speaker recognition
- **multi-modal** recognition
  - use anything

# why is this a **difficult task**?

- **missing** training data

35% of test speakers have no training data  
even worse for face recognition

- **asynchronous** data

*a speaker may be heard but not seen  
two (or more) people may appear at the same time  
people rarely pronounces their own name*

simple score fusion is not an option

- **heterogeneous** data

voice, face, text from OCR and from ASR

needs expertise in speech processing, computer vision and machine learning

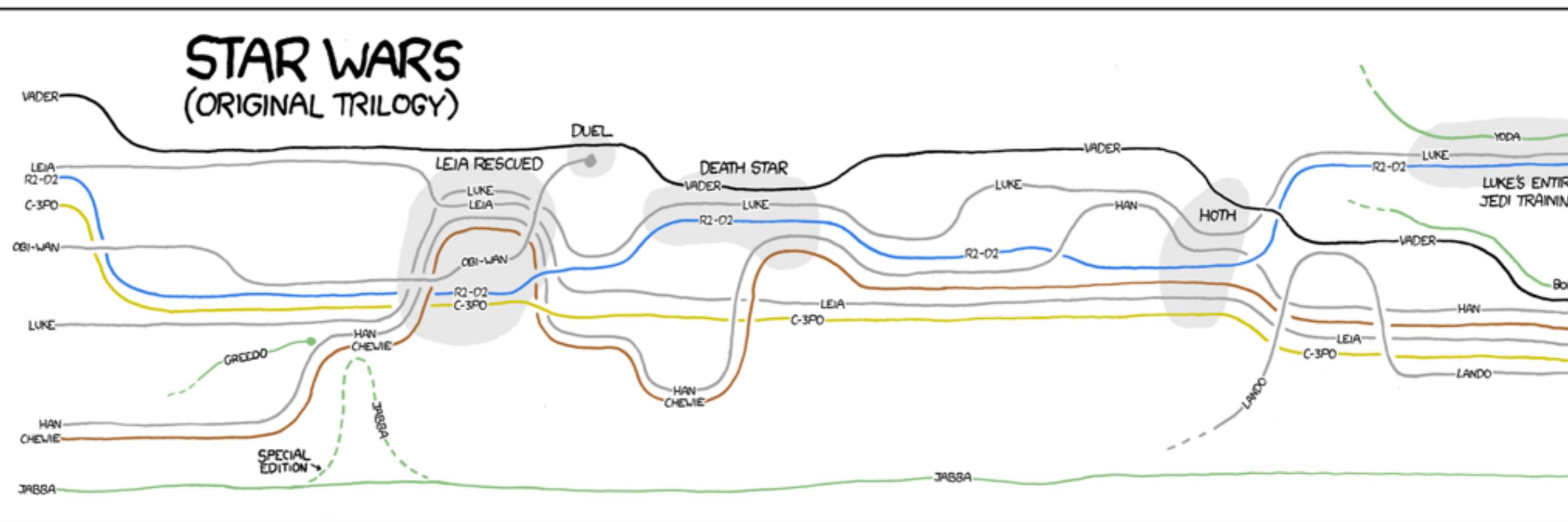
# why is this a useful task?

- **information retrieval**  
find all sequences where Sarkozy talks about immigration
- **radio/TV broadcast monitoring**  
measure total speech time for every politician (CSA)
- **second-screen** applications  
narrative graphs of TV/movie series



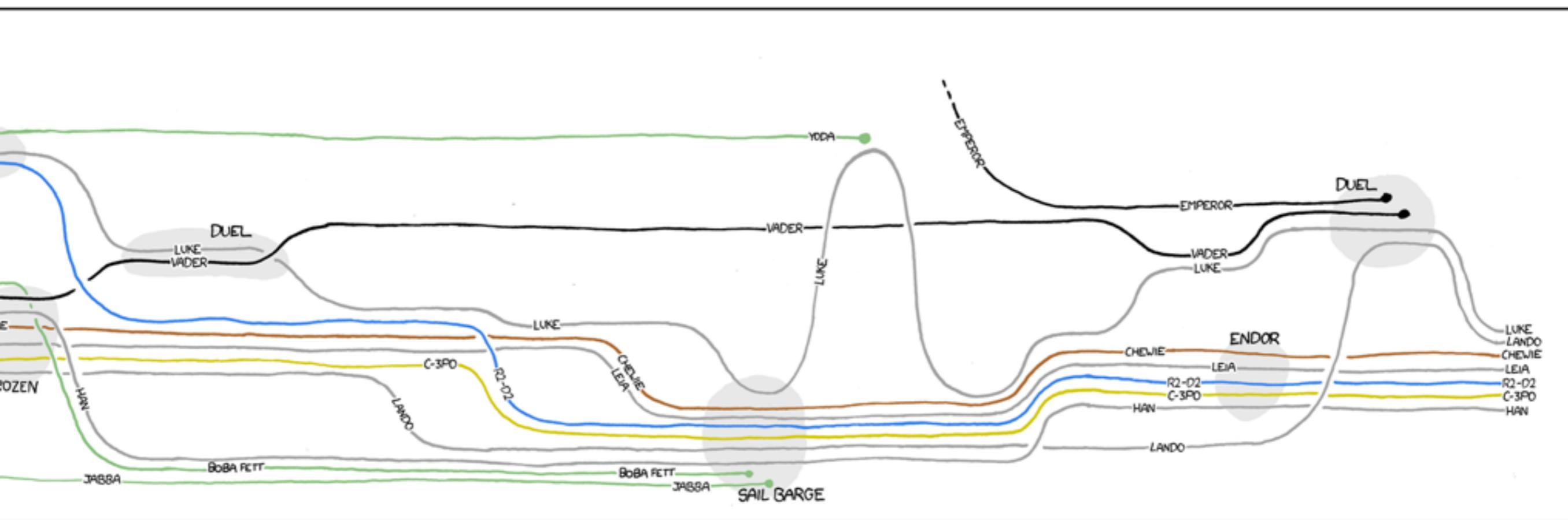
[xkcd.com/657/](http://xkcd.com/657/)

# narrative graphs of TV/movie series



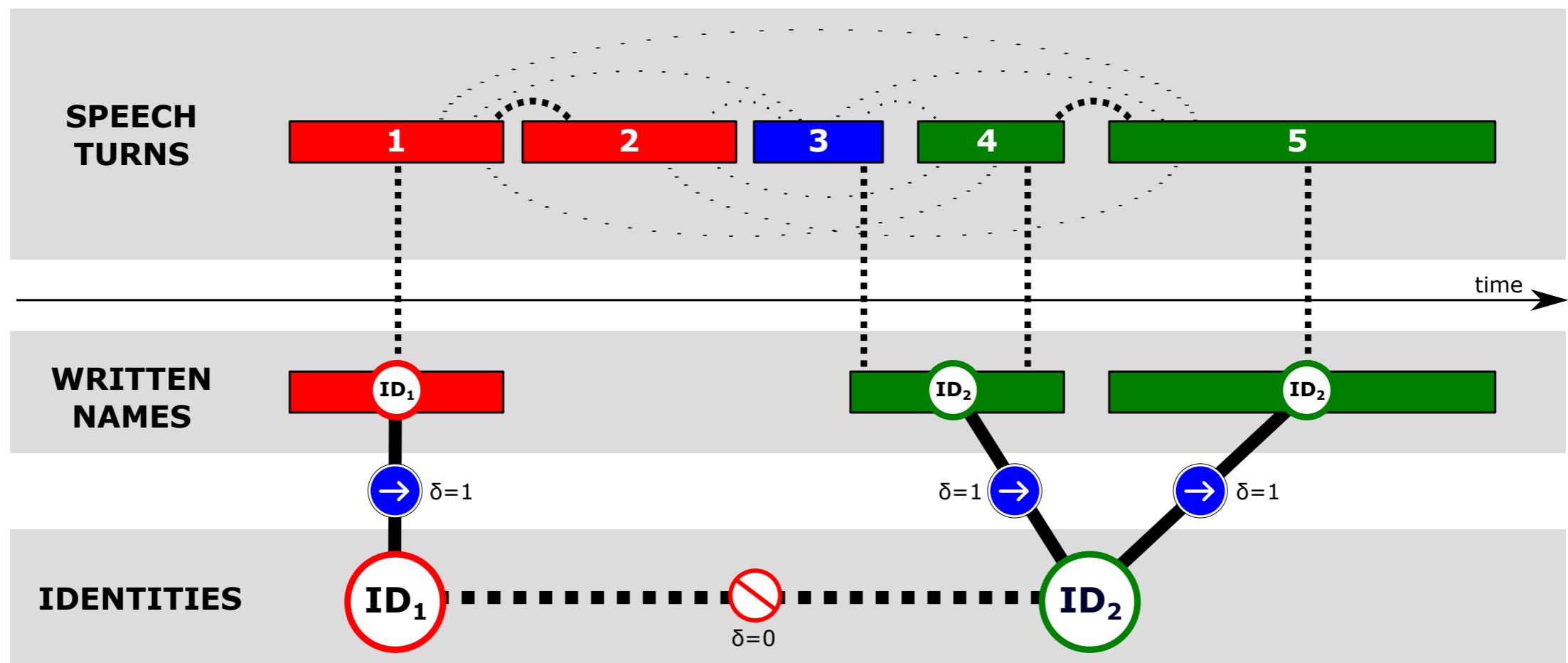
[xkcd.com/657/](http://xkcd.com/657/)

# narrative graphs of TV/movie series



[xkcd.com/657/](http://xkcd.com/657/)

# building person instance graphs



# person instance graph

- weighted undirected graph

$$\mathcal{G} = (\mathcal{V}, \mathcal{E}, p)$$

- instance vertices

speech turns  $t \in \mathcal{T}$

written names  $w \in \mathcal{W}$

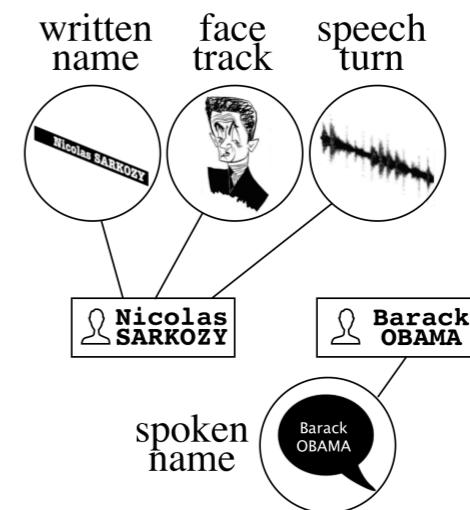
spoken names  $s \in \mathcal{S}$

~~face tracks~~

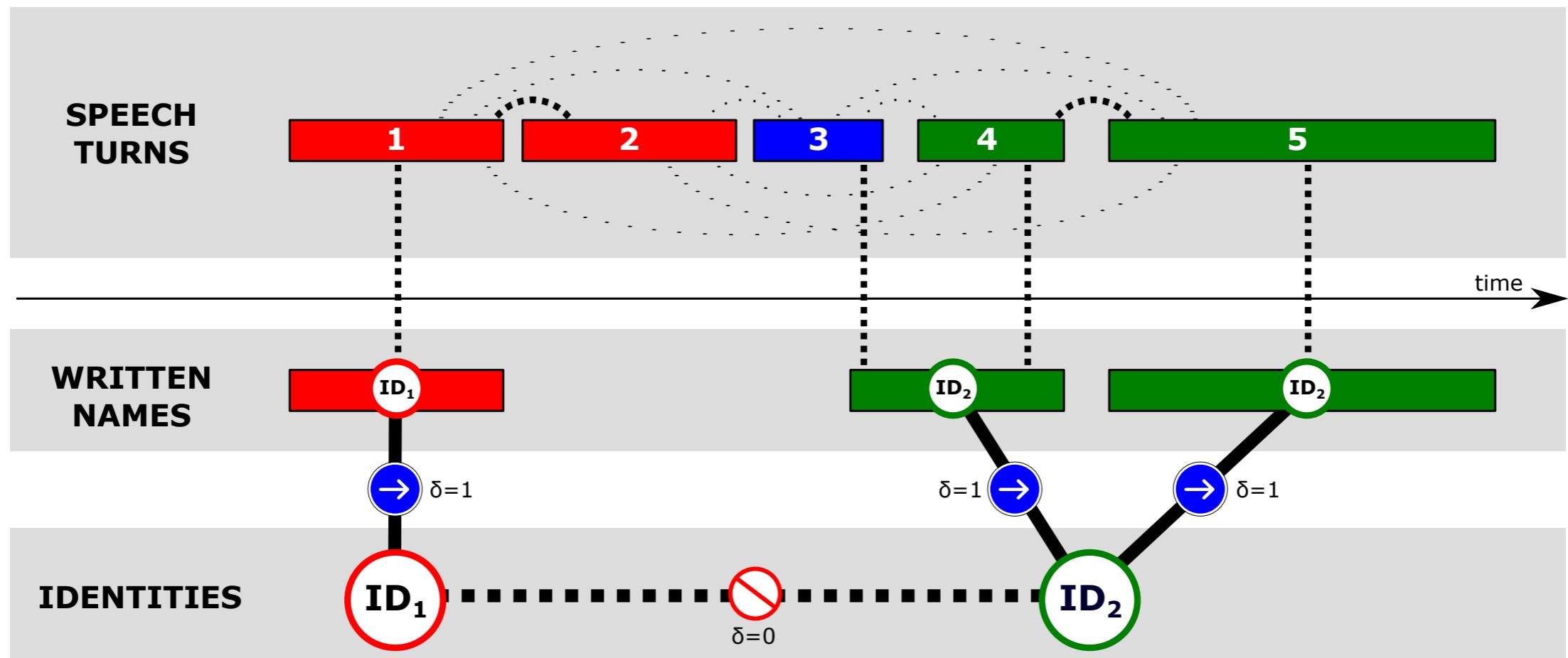
- identity vertices

- edges  $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$

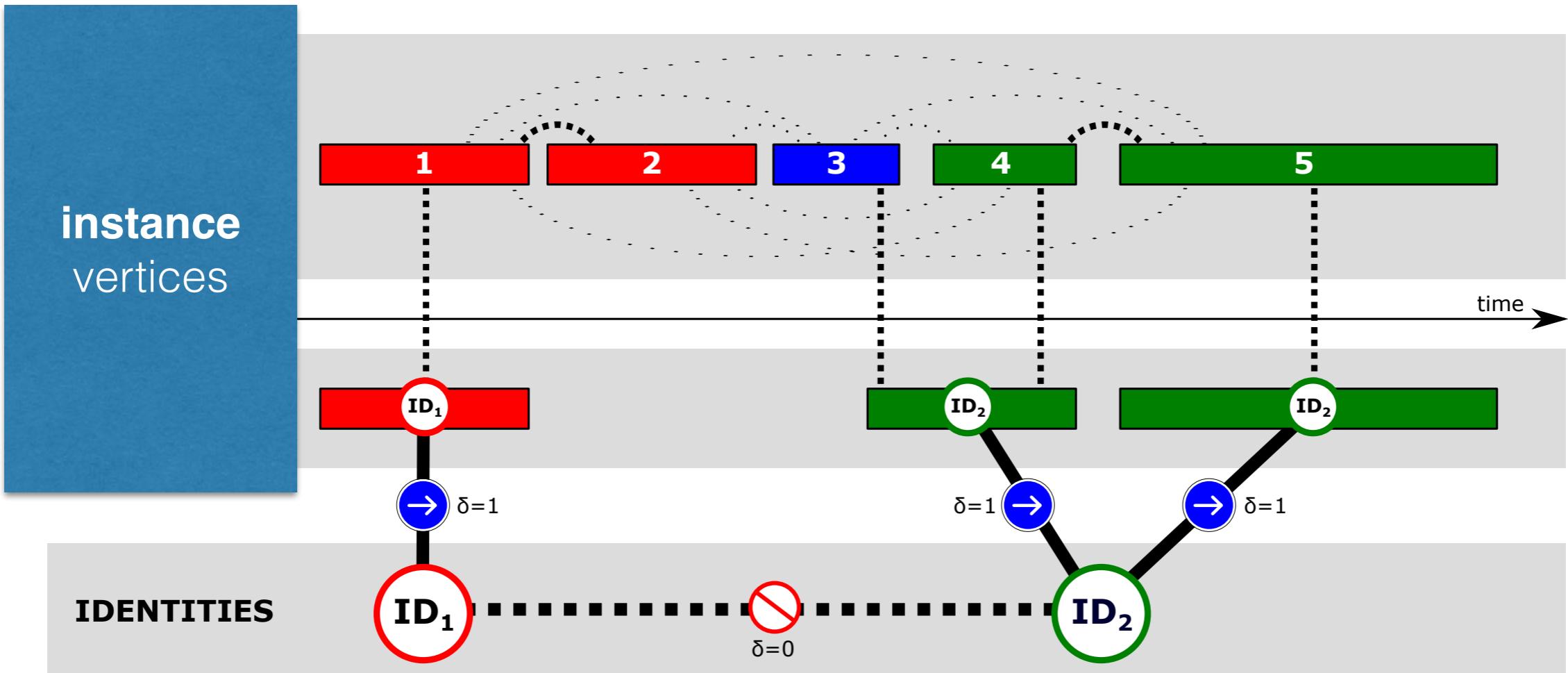
weighted by  $p \in [0, 1]^{\mathcal{E}}$



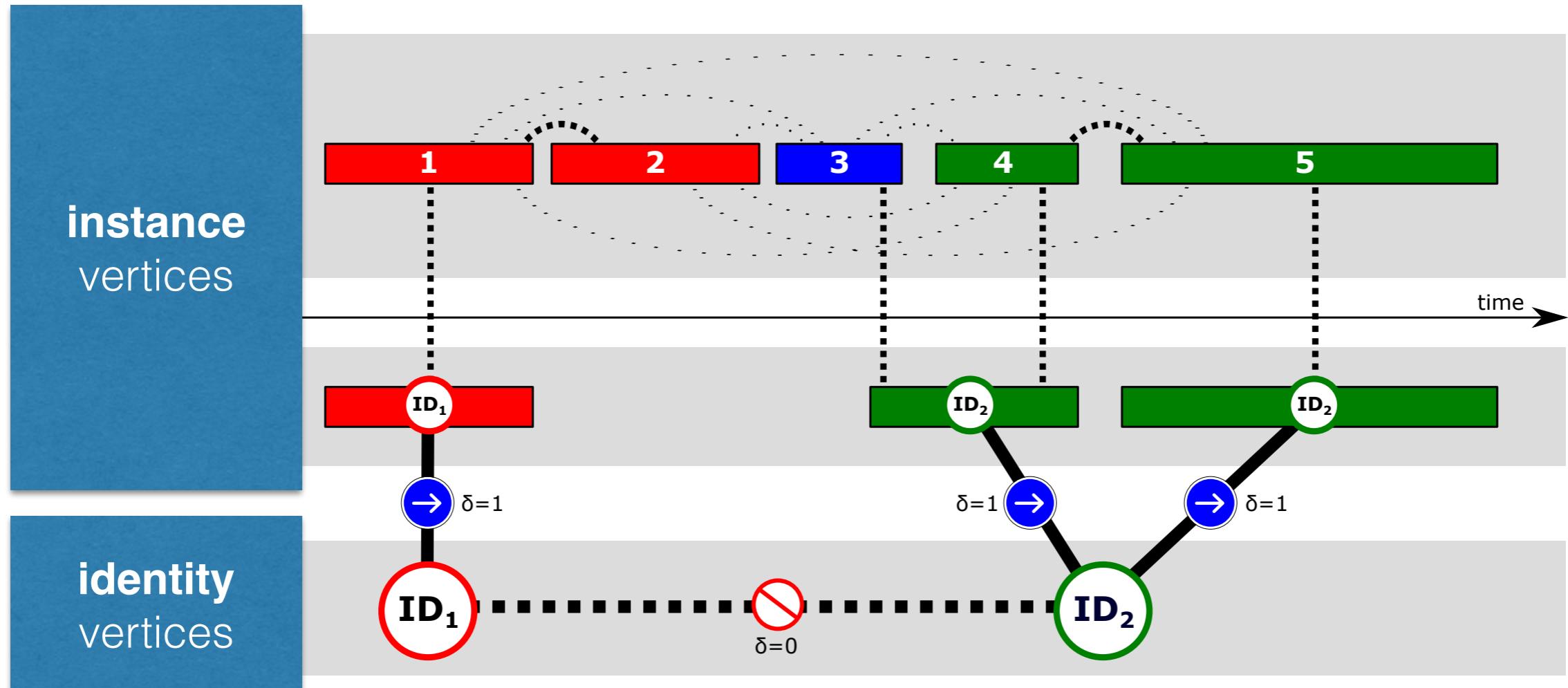
# vertices



# vertices

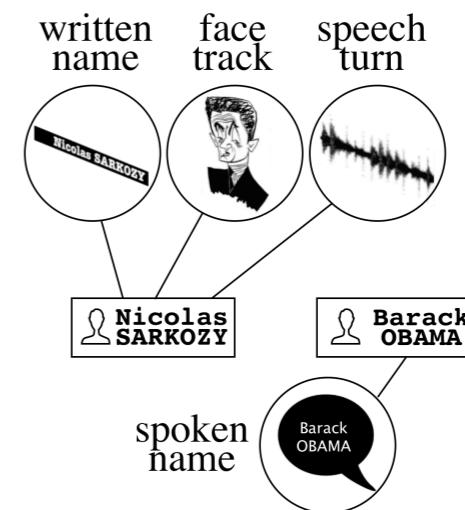


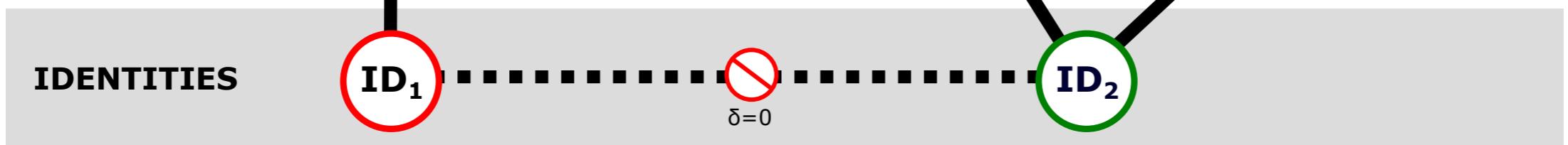
# vertices



# instance vertices

- **four potential sources of instance** vertices
  - **speech turns**  $t \in \mathcal{T}$   
speech activity detection  
divergence-based segmentation
  - **face tracks** (not used here)  
face detection and tracking
  - **written names**  $w \in \mathcal{W}$   
video optical character recognition  
named entity detection
  - **spoken names**  $s \in \mathcal{S}$   
manual speech transcription  
manual named entity detection





# identity vertices

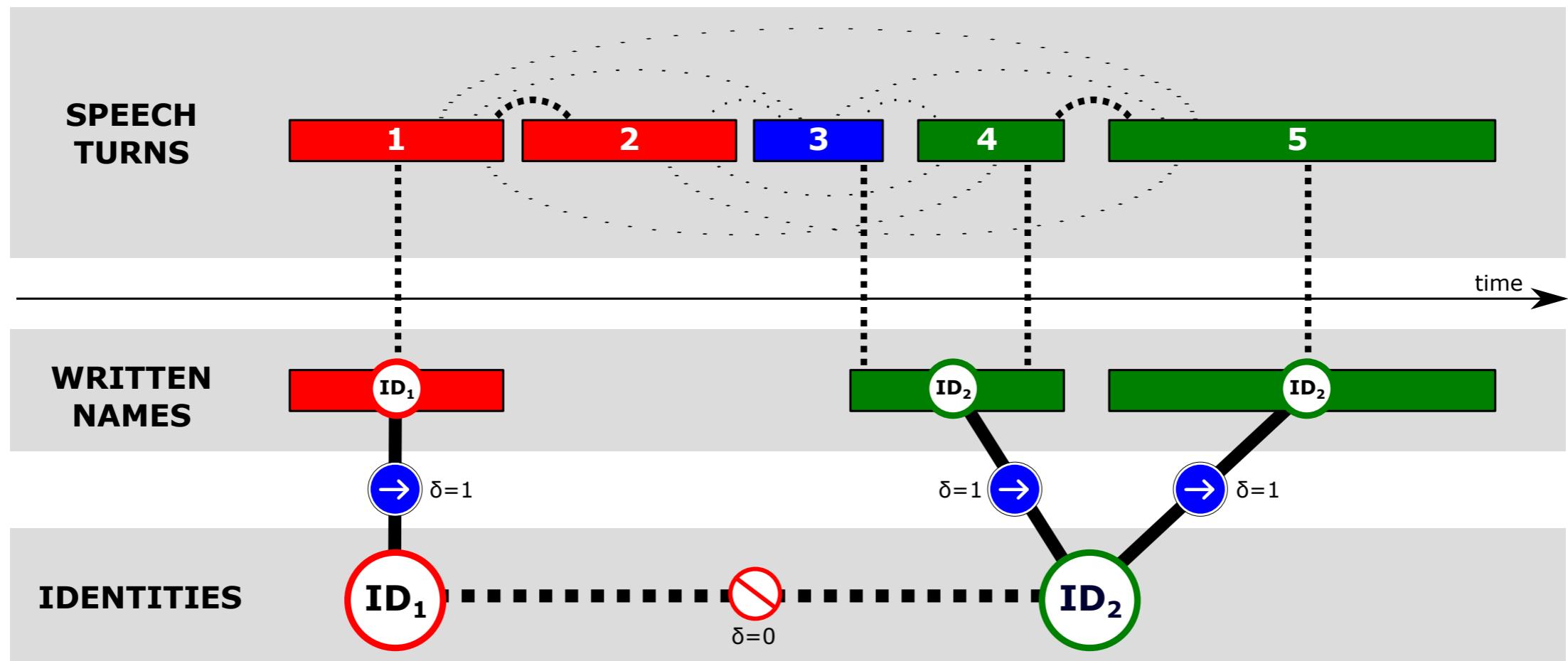
- **four potential sources of identity** vertices

**one unique identifier per person**

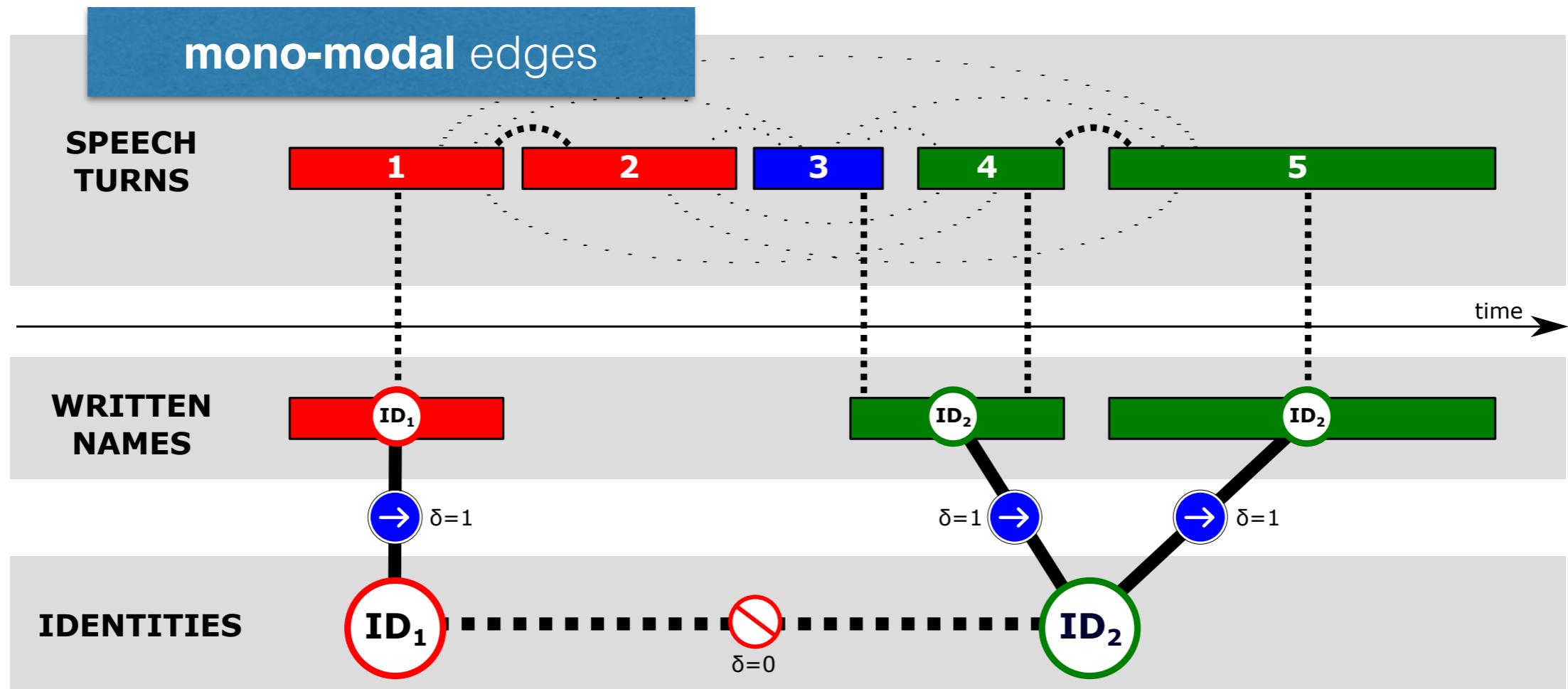
Nicolas Sarkozy  $\Rightarrow$  Nicolas\_SARKOZY

- **one identity vertex** per prior **speaker model**  $\mathcal{I}^*$   
350 speaker models  $\Rightarrow$  350 identity vertices
- **one identity vertex** per prior **face model**  
not used here
- **one identity vertex** for each **written name transcription**  $\mathcal{I}_{\mathcal{W}}$
- **one identity vertex** for each **spoken name transcription**  $\mathcal{I}_{\mathcal{S}}$
- $\mathcal{I} = \mathcal{I}_{\mathcal{W}} \cup \mathcal{I}_{\mathcal{S}} \cup \mathcal{I}^*$   
an identity vertex can belong to more than one subset...

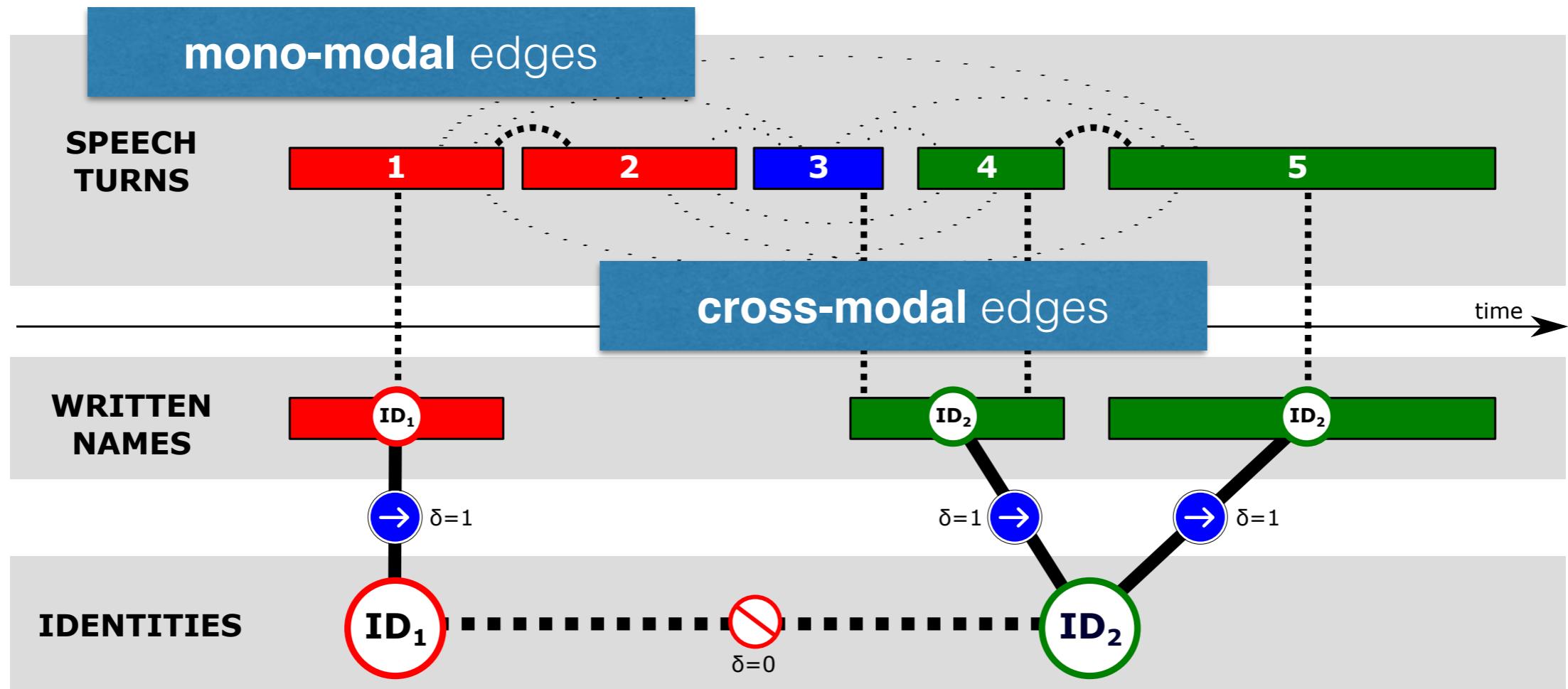
# edges



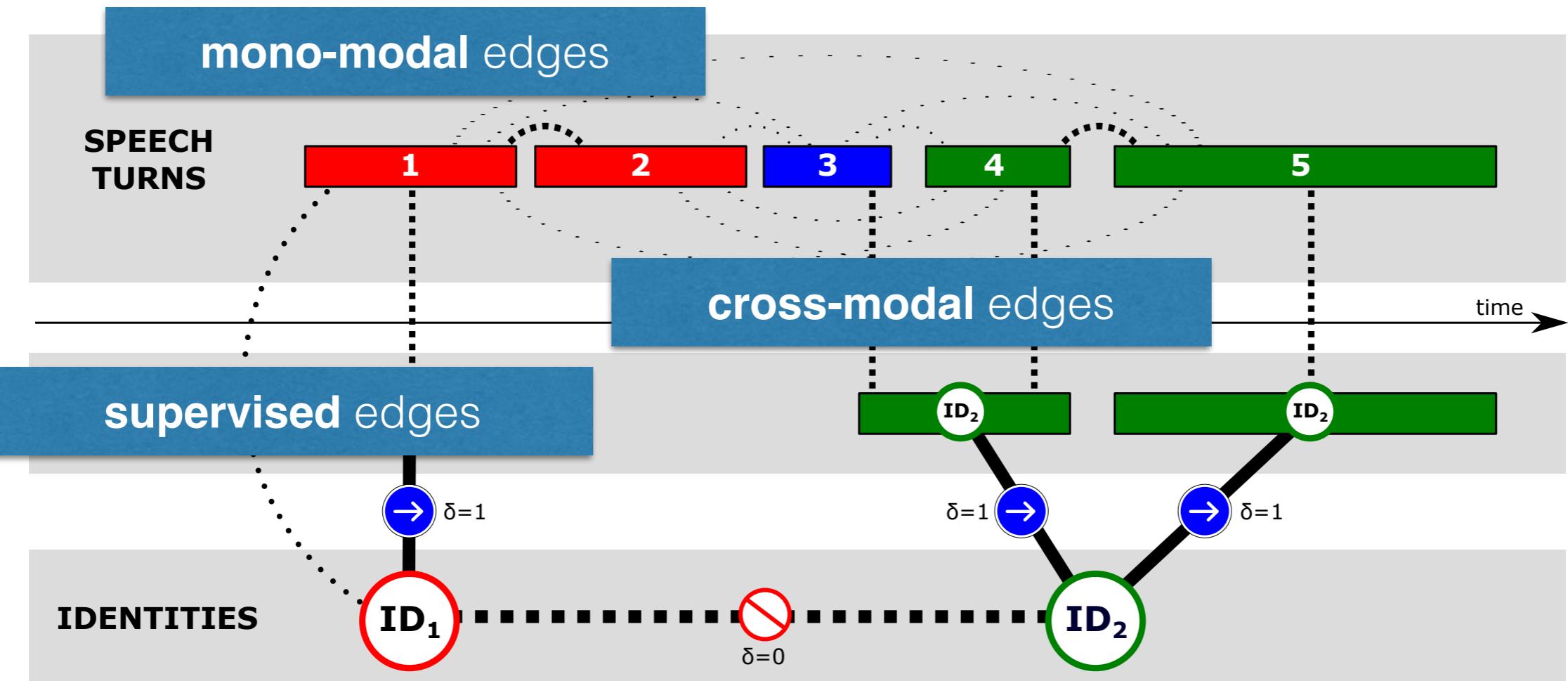
# edges



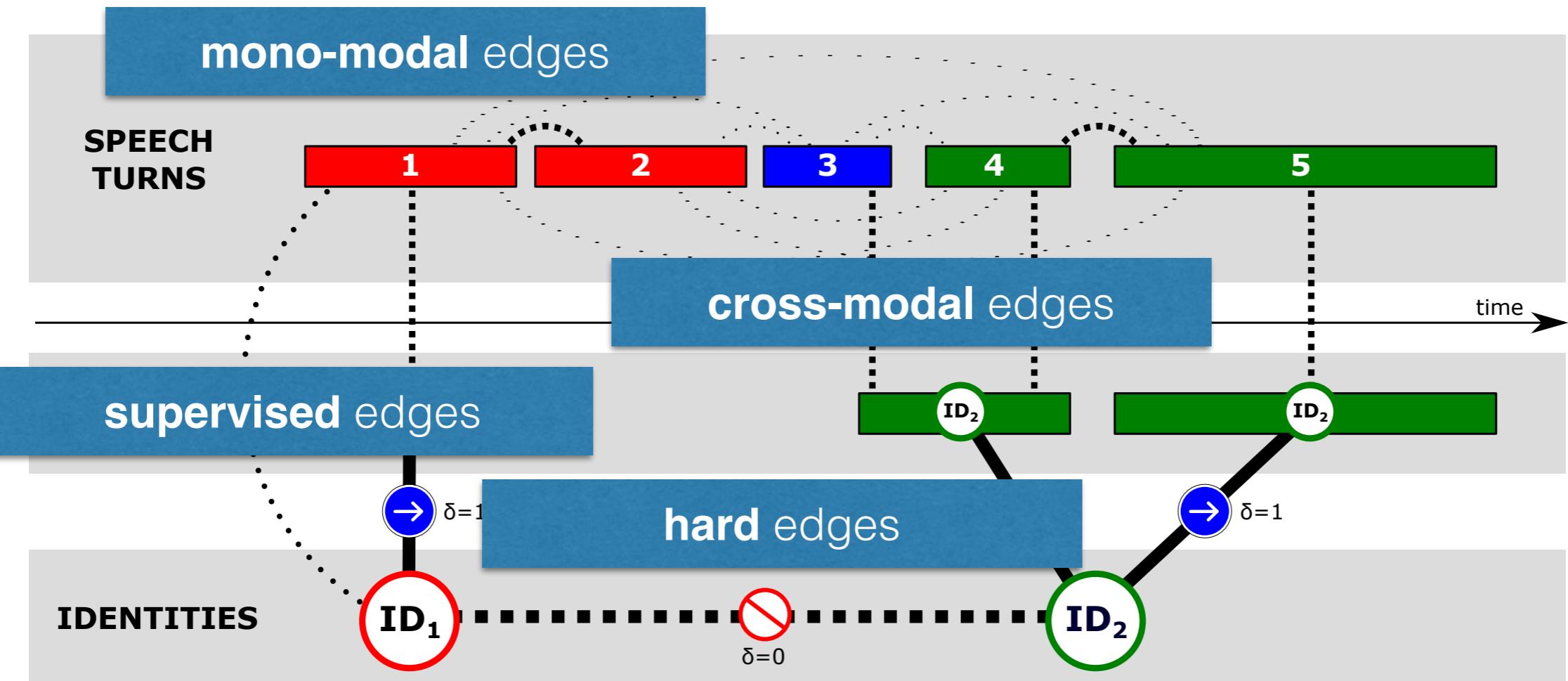
# edges



# edges



# edges

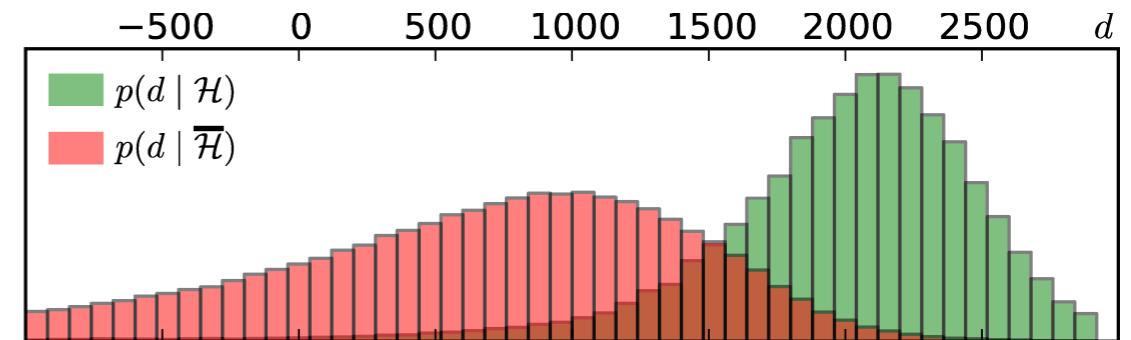


# mono-modal edges

(similarity between **speech turns**)

- Standard MFCC features
- Bayesian information criterion (BIC)

$$\begin{aligned} d_{tt'} = & (n_t + n_{t'}) \log |\Sigma_{t+t'}| \\ & - n_t \log |\Sigma_t| - n_{t'} \log |\Sigma_{t'}| \\ & - \frac{1}{2} \cdot \lambda \cdot \left( D + \frac{1}{2} D(D+1) \right) \log (n_t + n_{t'}) \end{aligned}$$

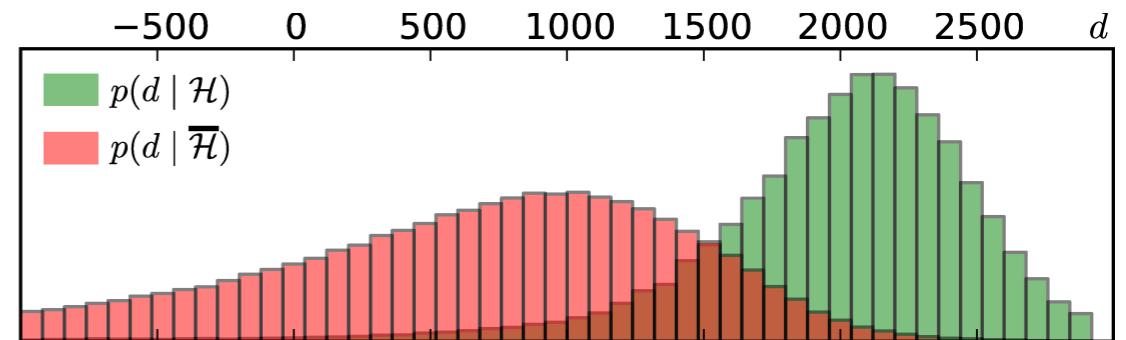


# mono-modal edges

(similarity between **speech turns**)

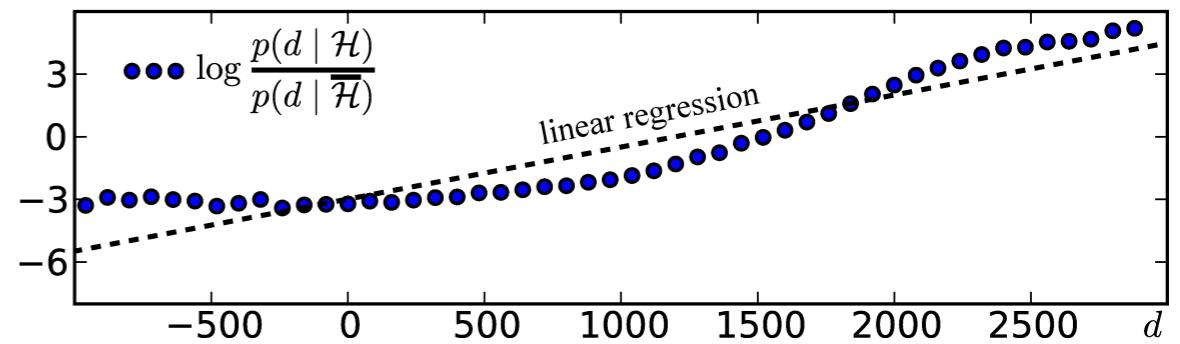
- Standard MFCC features
- Bayesian information criterion (BIC)

$$\begin{aligned}
 d_{tt'} &= (n_t + n_{t'}) \log |\Sigma_{t+t'}| \\
 &\quad - n_t \log |\Sigma_t| - n_{t'} \log |\Sigma_{t'}| \\
 &\quad - \frac{1}{2} \cdot \lambda \cdot \left( D + \frac{1}{2} D(D+1) \right) \log (n_t + n_{t'})
 \end{aligned}$$

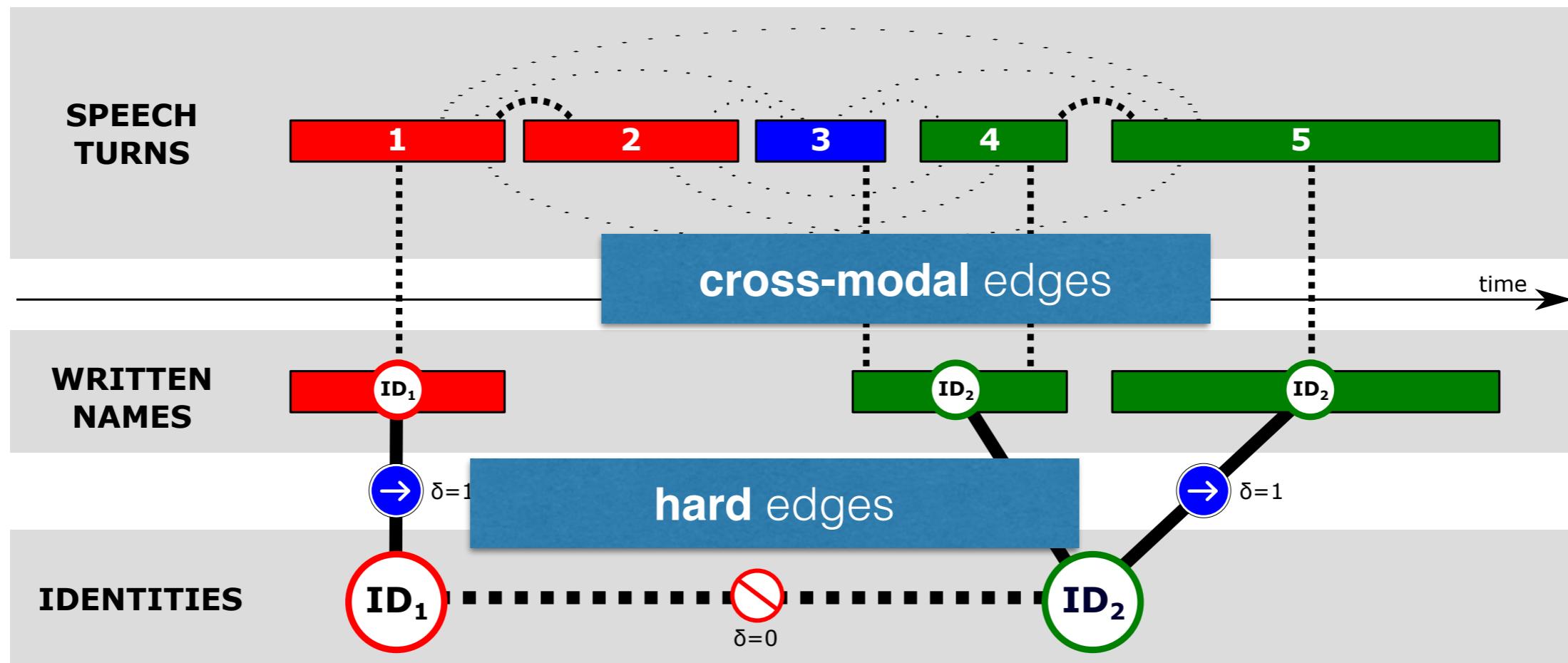


- Bayes' theorem
- Linear regression for log-likelihood ratio estimation

$$\begin{aligned}
 p_{tt'} &= p(\text{ID}(t) = \text{ID}(t') \mid d_{tt'}) \\
 &= \frac{1}{1 + \frac{\pi_{\neq}}{\pi_{=}} \cdot \frac{p(d_{tt'} \mid \text{ID}(t) \neq \text{ID}(t'))}{p(d_{tt'} \mid \text{ID}(t) = \text{ID}(t'))}}
 \end{aligned}$$



# edges



# cross-modal edges

$t \leftrightarrow w$

probability that a name written at time  $t$   
is the name of current speaker



$$p_{tw} = 0.956$$



$$p_{tw} = 0.996$$

# cross-modal edges

$t \leftrightarrow w$

probability that a name written at time t  
is the name of current speaker



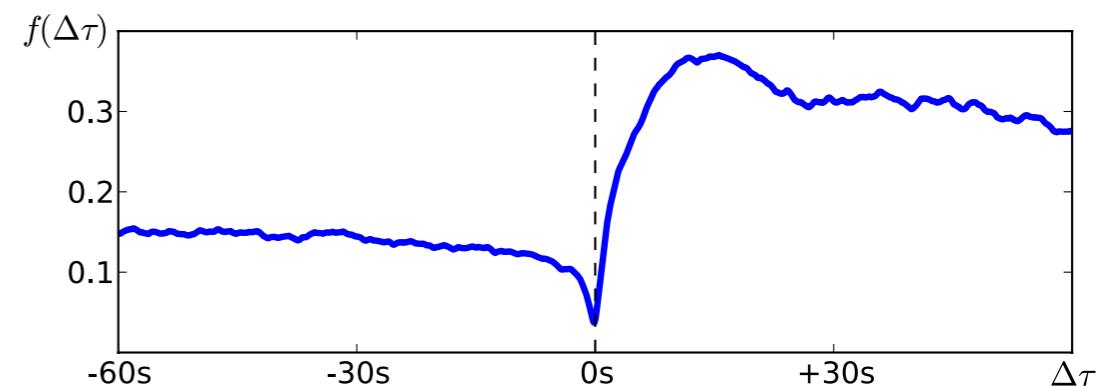
$$p_{tw} = 0.956$$



$$p_{tw} = 0.996$$

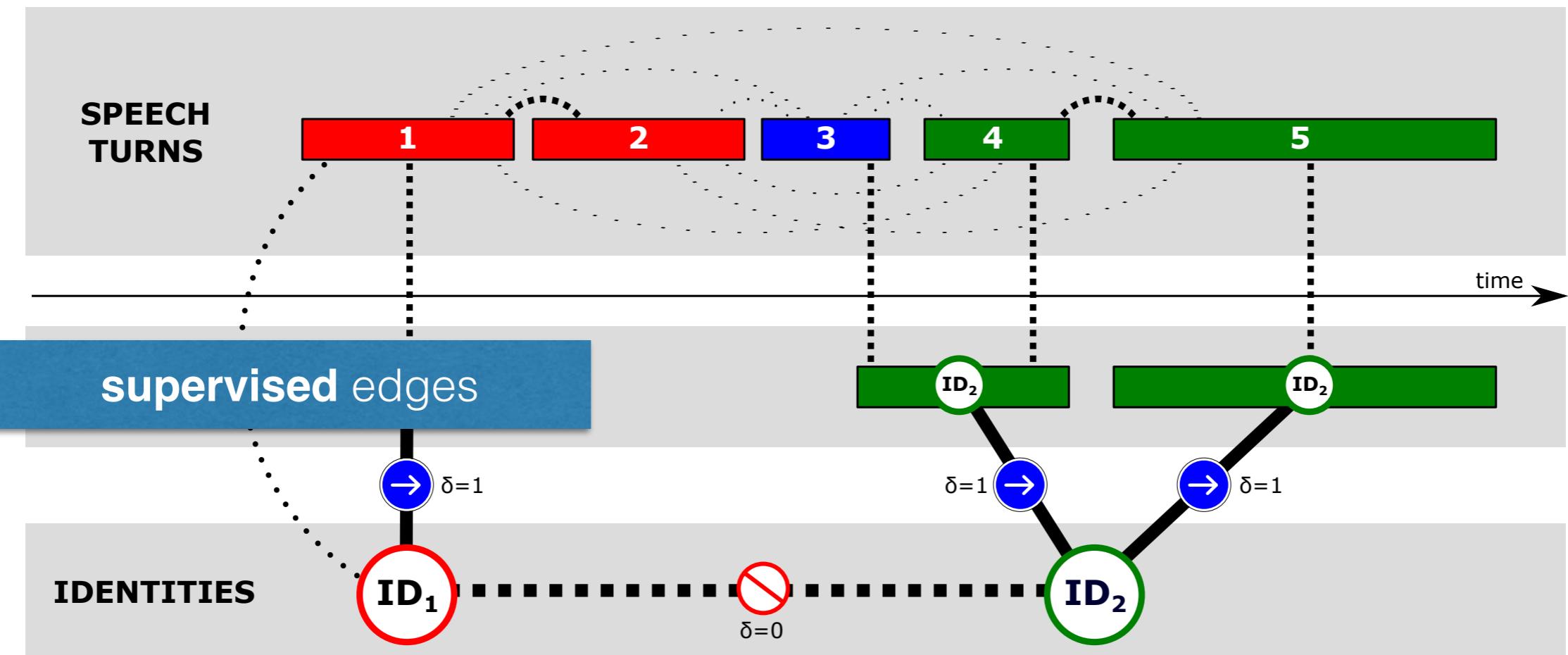
$t \leftrightarrow s$

probability that a name pronounced at time t  
is the name of speaker at time  $t + dt$



$$\begin{aligned} p_{ts} &= p(\text{ID}(t) = \text{ID}(s) \mid t, s) \\ &= \frac{1}{|T(t)|} \int_{\tau \in T(t)} f(\tau - \tau_s) \, d\tau \end{aligned}$$

# edges



# supervised edges

- GMM / UBM speaker identification

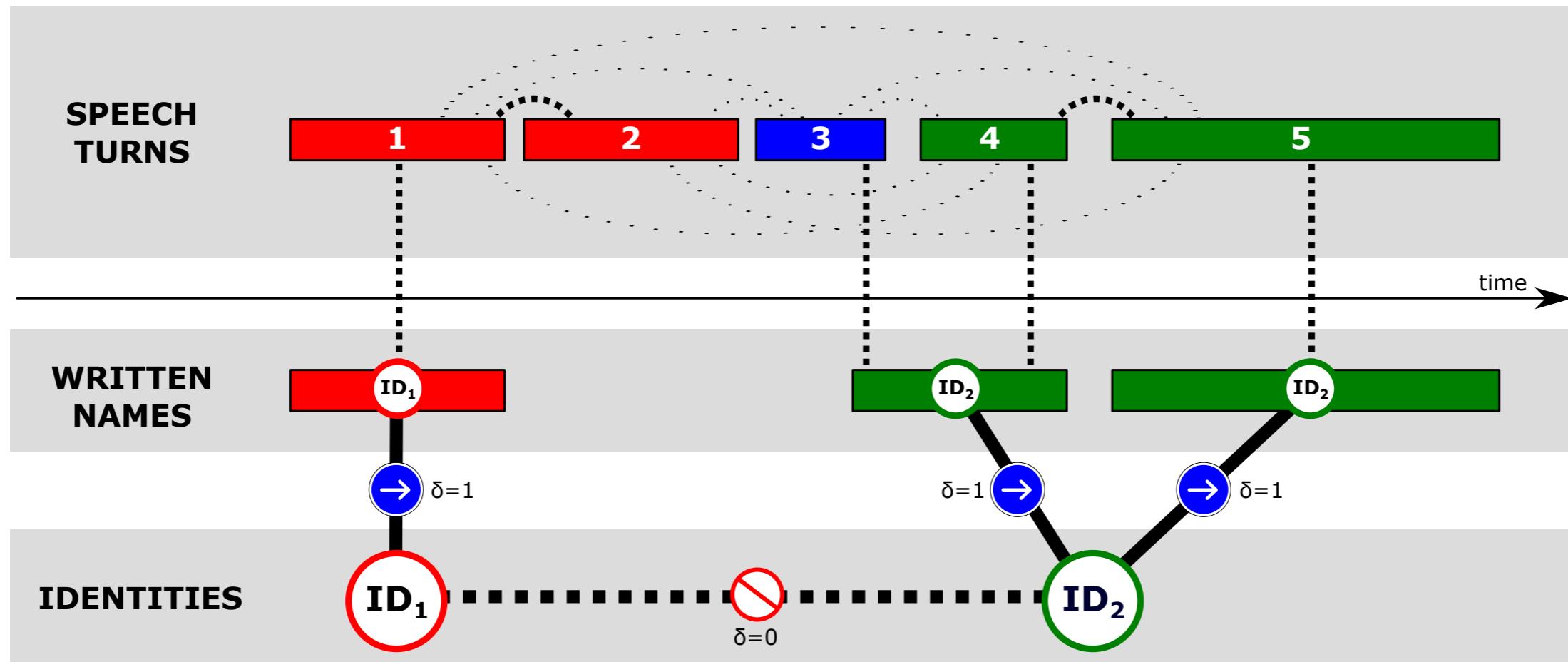
$$d_{ti} = \frac{1}{n_t} \left[ \log \prod_{x \in X_t} f(x | \lambda_i) - \log \prod_{x \in X_t} f(x | \lambda_\Omega) \right]$$

$$\begin{aligned} p_{ti} &= p(\text{ID}(t) = i \mid d_{ti}) \\ &= \frac{\pi_i \cdot \frac{p(d_{ti} \mid \text{ID}(t) = i)}{p(d_{ti} \mid \text{ID}(t) \neq i)}}{\pi_{\text{?}} + \sum_{i' \in \mathcal{I}^*} \pi_{i'} \cdot \frac{p(d_{ti'} \mid \text{ID}(t) = i')}{p(d_{ti'} \mid \text{ID}(t) \neq i')}} \end{aligned}$$

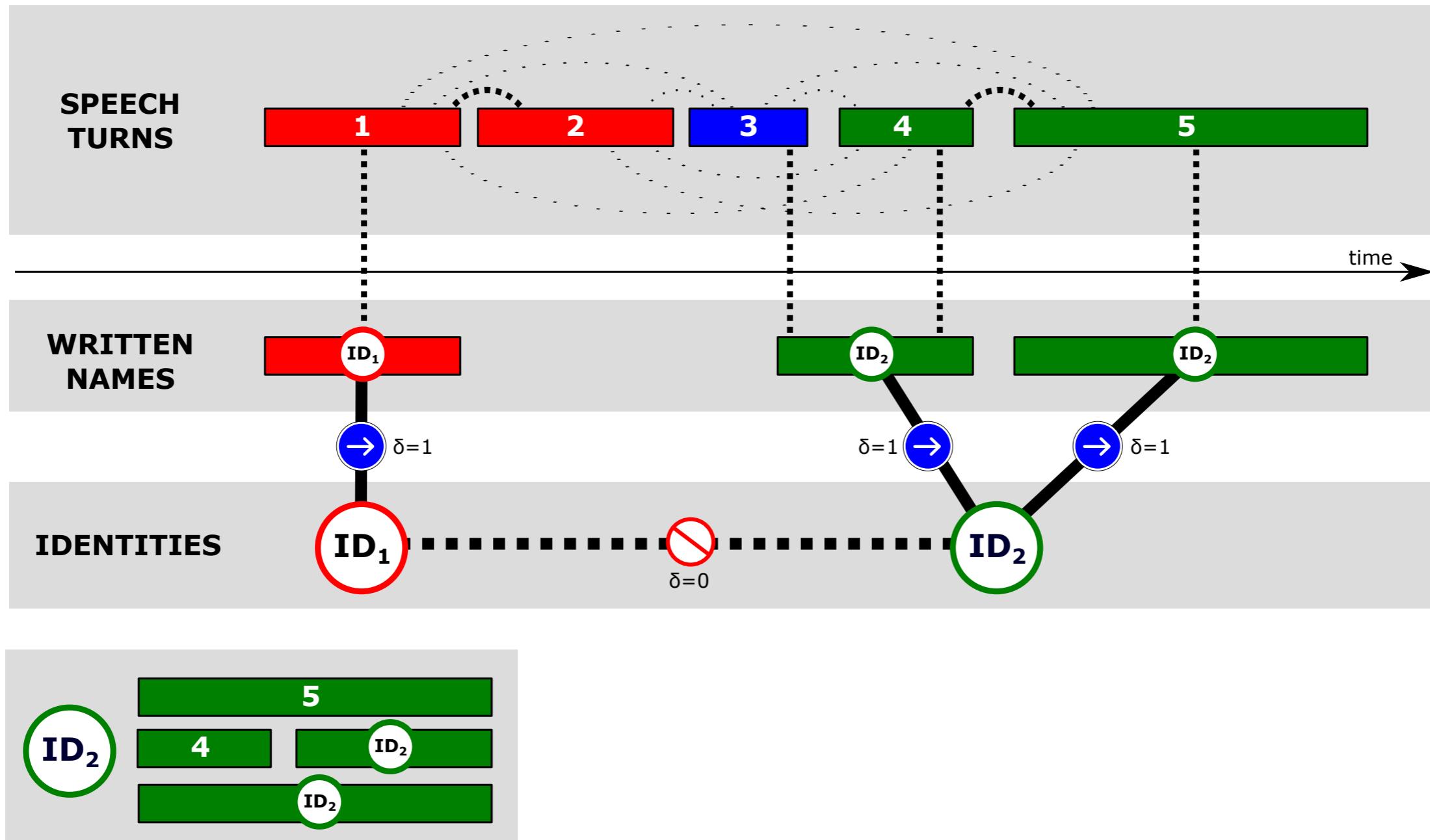
- every speech turn is connected  $t \leftrightarrow i$  to all ( $\approx 350$ ) identity vertices

# mining person instance graphs

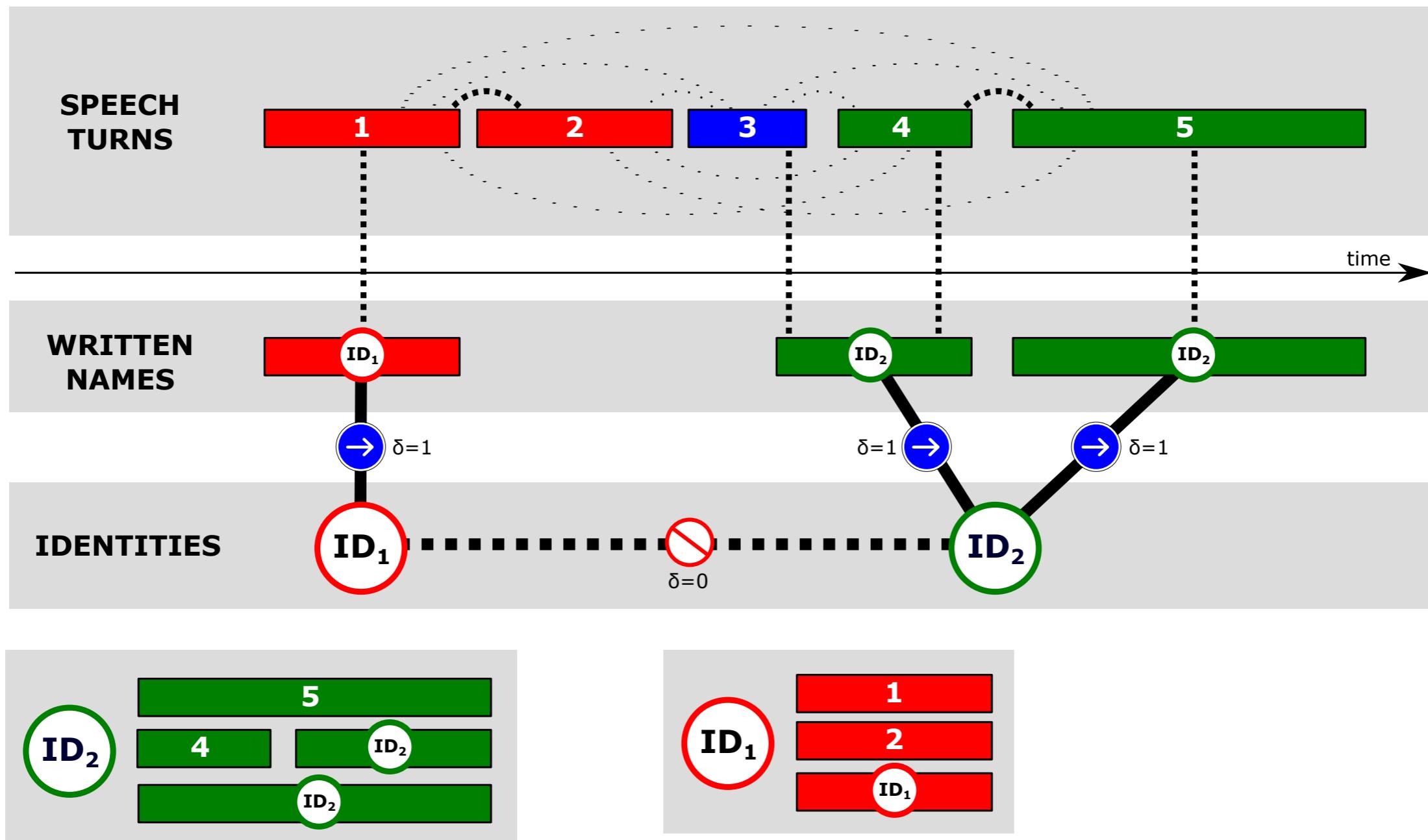
# expected results



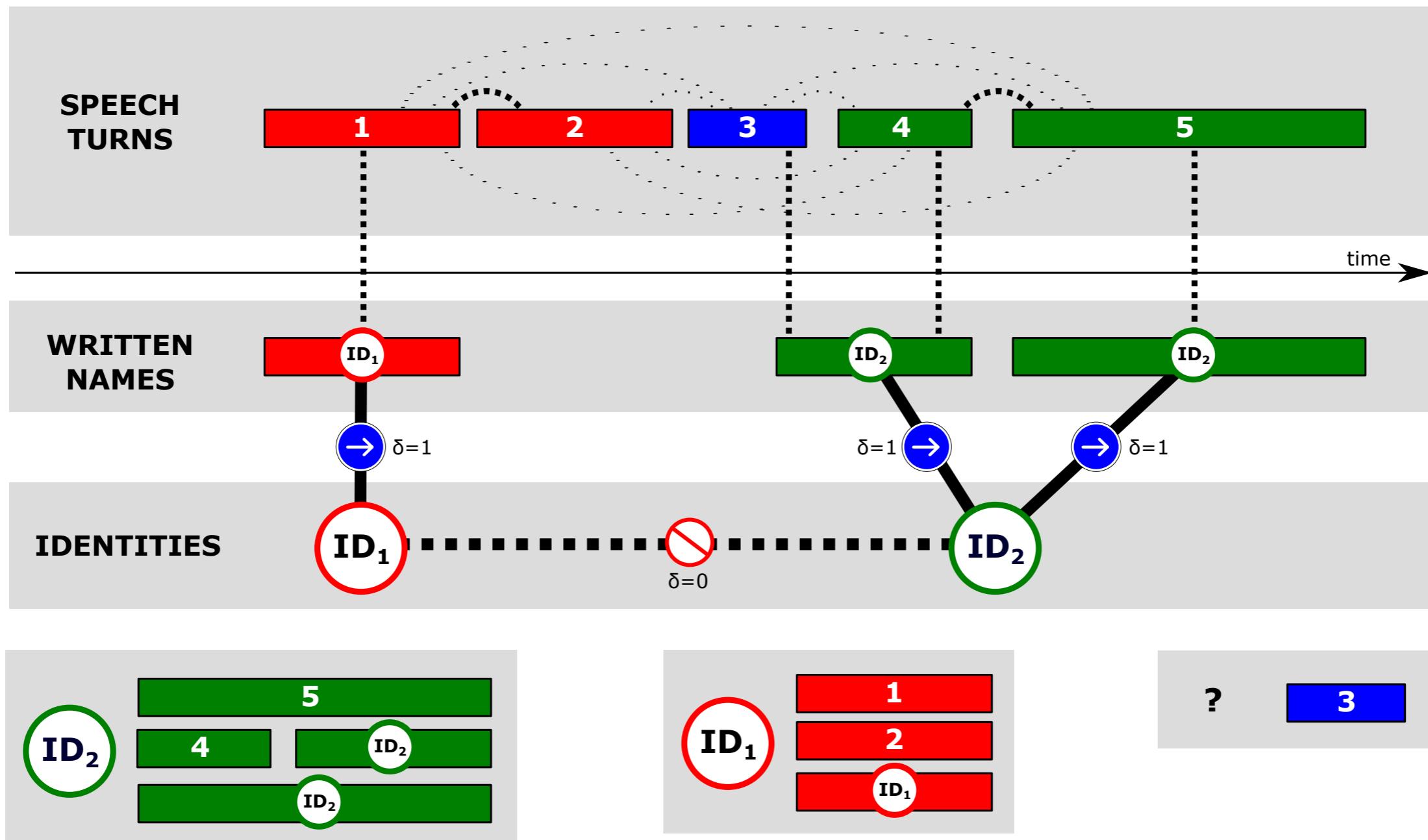
# expected results



# expected results

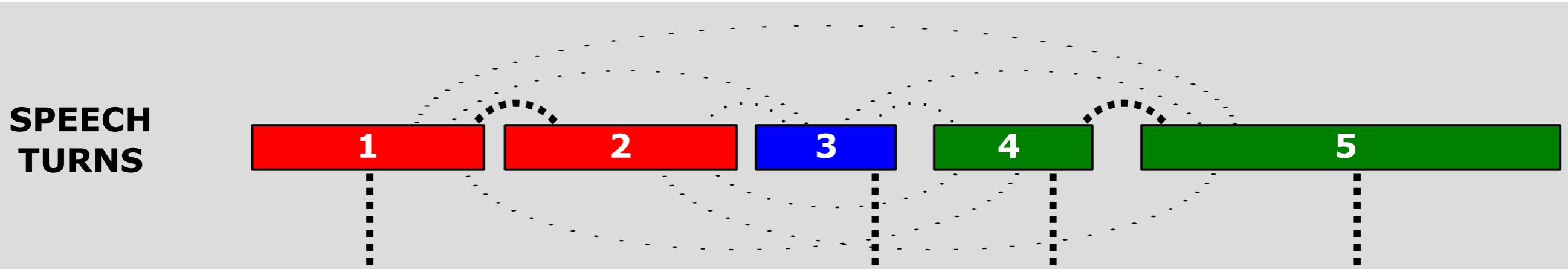


# expected results



# speaker diarization

- the task of **partitioning** and **labeling** an audio stream into homogeneous **speech segments** according to the **identity of the speaker**
- usually addressed as a **segmentation** problem followed by a **clustering** problem



# three limitations

- k-means (et al.) rely on the assumption that the **number of clusters** is known ***a priori***.
- hierarchical agglomerative clustering does not guarantee **global optimality** does not provide clear **stopping criterion**
- cannot deal with **incomplete similarity matrices**
- *Finkel & Manning* formulate clustering as an **integer linear programming** problem (for co-reference resolution)

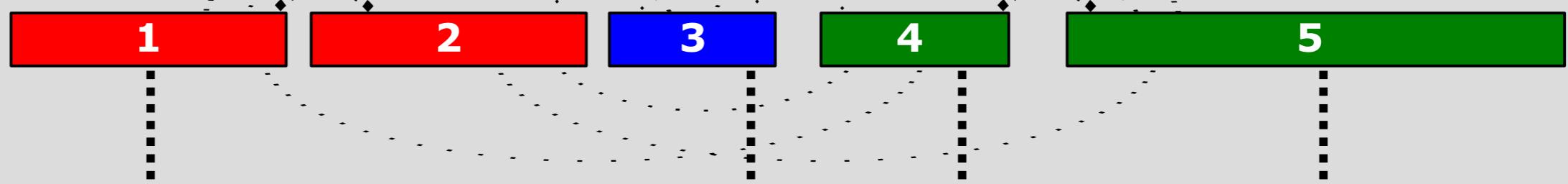
# notations

- $\mathcal{G} = (\mathcal{V}, \mathcal{E}, p)$  is an **undirected weighted** graph
- $\mathcal{V}$  is the set of **vertices (one per item** to cluster)
- $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$  is the set of **edges**  
two vertices **may or may not be connected**
- edges are weighted by the **probability** that  
connected vertices belong to the **same cluster**

$$p: \mathcal{E} \rightarrow [0, 1]$$

$$(v, v') \mapsto p_{vv'} = p(\text{ID}(v) = \text{ID}(v') \mid v, v')$$

SPEECH  
TURNS



# clustering function

- Any output of a valid clustering algorithm can be described by a clustering function

$$\delta: \mathcal{V} \times \mathcal{V} \rightarrow \{0, 1\}$$

$$(v, v') \mapsto \begin{cases} 1 & \text{if } v \text{ and } v' \text{ are in the same cluster,} \\ 0 & \text{otherwise.} \end{cases}$$

# clustering function

- Any output of a valid clustering algorithm can be described by a clustering function

$$\delta: \mathcal{V} \times \mathcal{V} \rightarrow \{0, 1\}$$

$$(v, v') \mapsto \begin{cases} 1 & \text{if } v \text{ and } v' \text{ are in the same cluster,} \\ 0 & \text{otherwise.} \end{cases}$$

- Reciprocally, additional constraints are needed in order to guarantee a valid clustering

$$\Delta_{\mathcal{V}} = \left\{ \begin{array}{ll} \delta \in \{0, 1\}^{\mathcal{V} \times \mathcal{V}} & \text{s.t. } \forall (v, v', v'') \in \mathcal{V}^3, \\ \text{(a) } \delta_{vv} = 1 & \text{reflexivity} \\ \text{(b) } \delta_{vv'} = \delta_{v'v} & \text{symmetry} \\ \text{(c) } \delta_{vv'} = 1 \wedge \delta_{v'v''} = 1 \implies \delta_{vv''} = 1 & \text{transitivity} \end{array} \right.$$

# objective function

- Find the **clustering function** that maximizes the objective function:  $\delta^* = \underset{\delta \in \Delta_V}{\operatorname{argmax}} \mathcal{L}^\alpha(\delta, \mathcal{E}, p)$

# objective function

- Find the **clustering function** that maximizes the objective function:  $\delta^* = \underset{\delta \in \Delta_{\mathcal{V}}}{\operatorname{argmax}} \mathcal{L}^\alpha(\delta, \mathcal{E}, p)$

$$\mathcal{L}^\alpha(\delta, \mathcal{E}, p) = |\mathcal{E}|^{-1} \left[ \alpha \cdot \overbrace{\sum_{(v,v') \in \mathcal{E}} \delta_{vv'} \cdot p_{vv'}}^{\text{intra-cluster similarity}} + (1 - \alpha) \cdot \overbrace{\sum_{(v,v') \in \mathcal{E}} (1 - \delta_{vv'}) \cdot (1 - p_{vv'})}^{\text{inter-cluster dissimilarity}} \right]$$

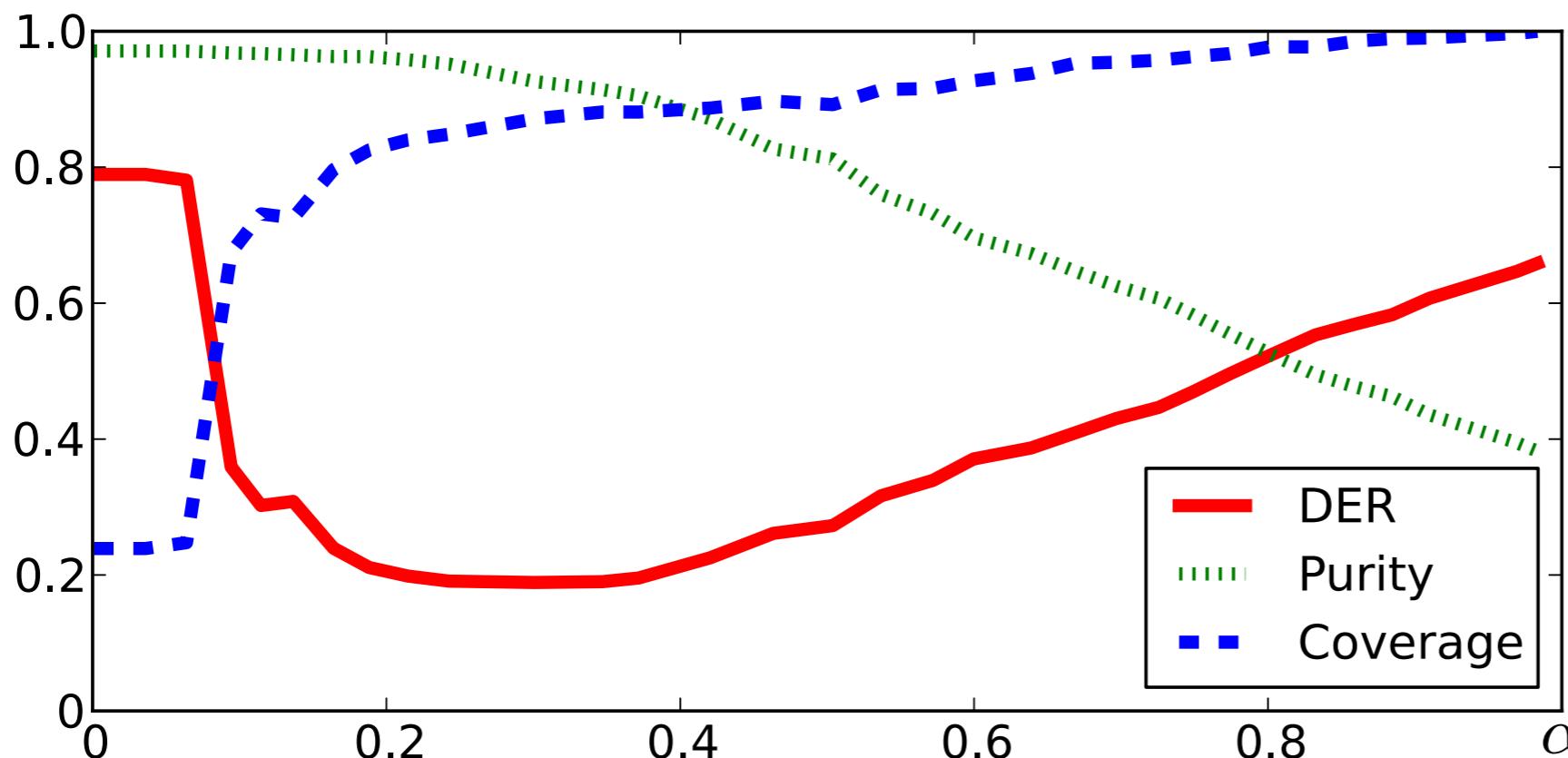
# effect of $\alpha$

$$\mathcal{L}^\alpha(\delta, \mathcal{E}, p) = |\mathcal{E}|^{-1} \left[ \underbrace{\alpha \cdot \sum_{(v,v') \in \mathcal{E}} \delta_{vv'} \cdot p_{vv'}}_{\text{intra-cluster similarity}} + (1 - \alpha) \cdot \underbrace{\sum_{(v,v') \in \mathcal{E}} (1 - \delta_{vv'}) \cdot (1 - p_{vv'})}_{\text{inter-cluster dissimilarity}} \right]$$

# effect of $\alpha$

$$\mathcal{L}^\alpha(\delta, \mathcal{E}, p) = |\mathcal{E}|^{-1} \left[ \underbrace{\alpha \cdot \sum_{(v,v') \in \mathcal{E}} \delta_{vv'} \cdot p_{vv'}}_{\text{intra-cluster similarity}} + (1 - \alpha) \cdot \underbrace{\sum_{(v,v') \in \mathcal{E}} (1 - \delta_{vv'}) \cdot (1 - p_{vv'})}_{\text{inter-cluster dissimilarity}} \right]$$

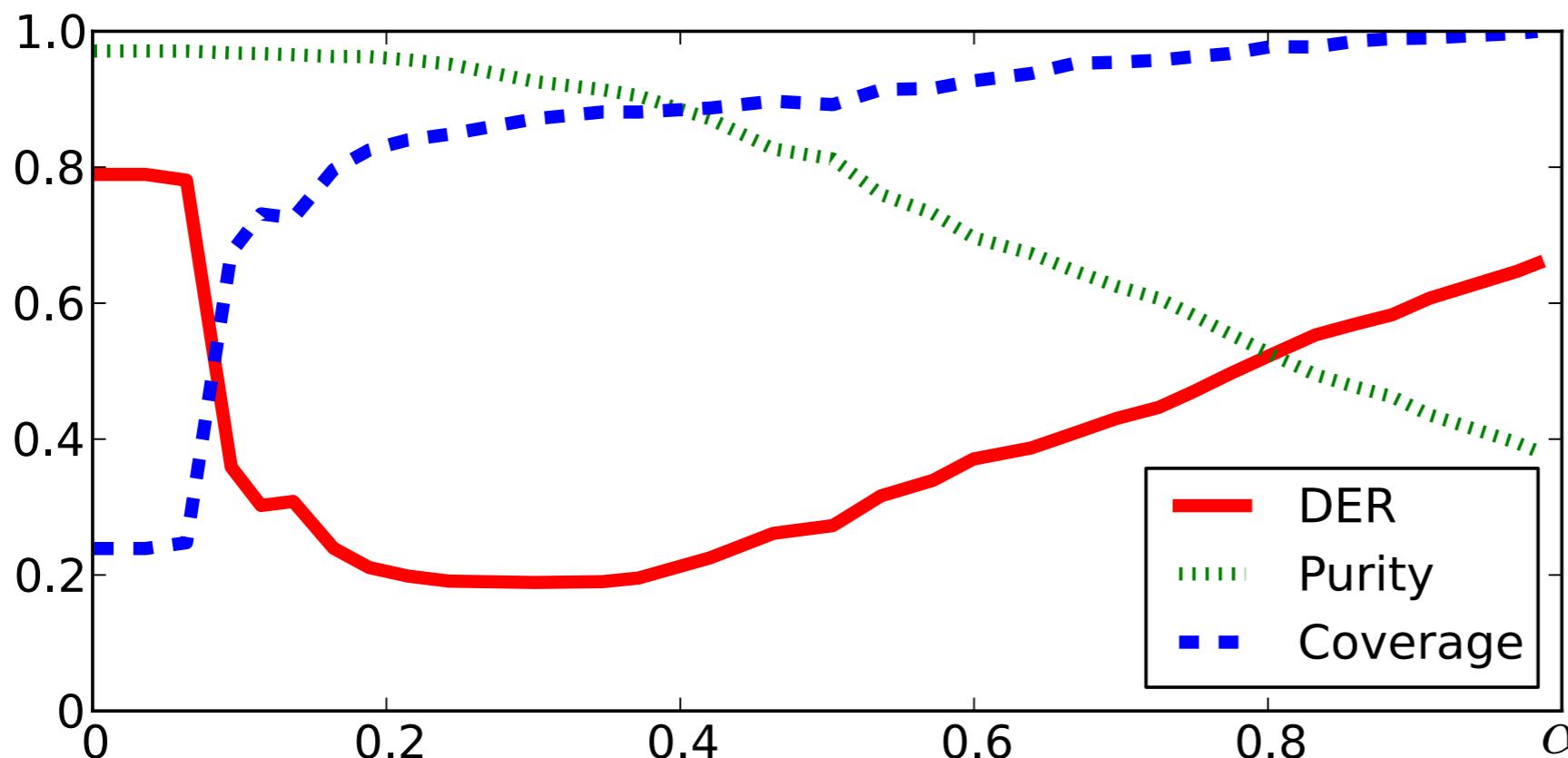
- cluster **purity** vs. cluster **coverage**



# effect of $\alpha$

$$\mathcal{L}^\alpha(\delta, \mathcal{E}, p) = |\mathcal{E}|^{-1} \left[ \underbrace{\alpha \cdot \sum_{(v,v') \in \mathcal{E}} \delta_{vv'} \cdot p_{vv'}}_{\text{intra-cluster similarity}} + \underbrace{(1 - \alpha) \cdot \sum_{(v,v') \in \mathcal{E}} (1 - \delta_{vv'}) \cdot (1 - p_{vv'})}_{\text{inter-cluster dissimilarity}} \right]$$

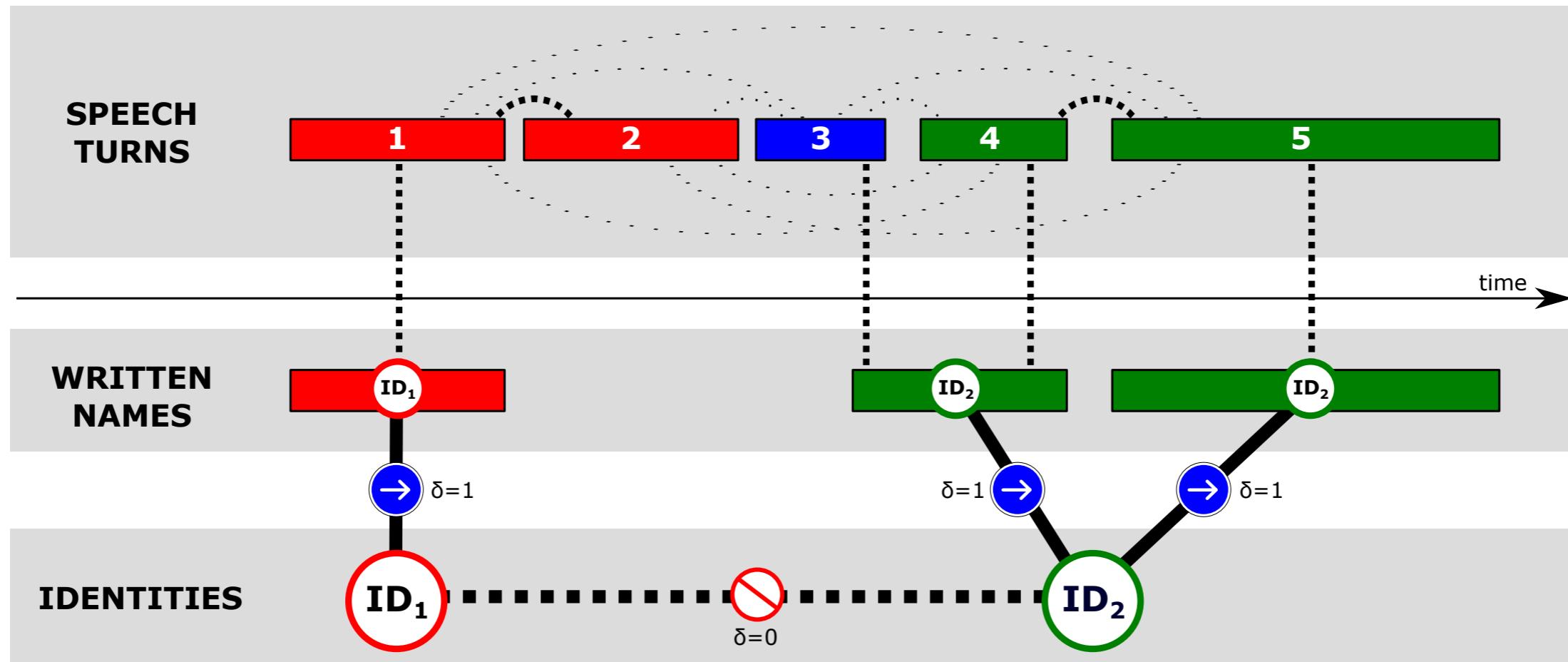
- cluster **purity** vs. cluster **coverage**



not nearly as good as standard BIC clustering



# extension to multiple modalities



# multimodal objective function

$$\mathcal{L}^\alpha(\delta, \mathcal{E}, p) = |\mathcal{E}|^{-1} \left[ \underbrace{\alpha \cdot \sum_{(v,v') \in \mathcal{E}} \delta_{vv'} \cdot p_{vv'}}_{\text{intra-cluster similarity}} + (1 - \alpha) \cdot \underbrace{\sum_{(v,v') \in \mathcal{E}} (1 - \delta_{vv'}) \cdot (1 - p_{vv'})}_{\text{inter-cluster dissimilarity}} \right]$$

would result in giving **more weight** to **dominant type** of edges

# multimodal objective function

$$\mathcal{L}^\alpha(\delta, \mathcal{E}, p) = |\mathcal{E}|^{-1} \left[ \underbrace{\alpha \cdot \sum_{(v,v') \in \mathcal{E}} \delta_{vv'} \cdot p_{vv'}}_{\text{intra-cluster similarity}} + (1 - \alpha) \cdot \underbrace{\sum_{(v,v') \in \mathcal{E}} (1 - \delta_{vv'}) \cdot (1 - p_{vv'})}_{\text{inter-cluster dissimilarity}} \right]$$

would result in giving **more weight** to **dominant type** of edges

$$\mathcal{L}_\beta^\alpha(\delta, \mathcal{E}, p) = \sum_{\substack{x \in \{\mathcal{T}, \mathcal{W}, \mathcal{S}, \mathcal{I}\} \\ y \in \{\mathcal{T}, \mathcal{W}, \mathcal{S}, \mathcal{I}\}}} \beta_{xy} \cdot \mathcal{L}^{\alpha_{xy}}(\delta, \mathcal{E} \cap (x \times y), p) \quad \alpha_{xy} \in [0, 1] \quad \beta_{xy} \in [0, 1] \quad \sum_{x,y} \beta_{xy} = 1$$

**one objective function per edge type subgraph**

allows different weights beta for different edge types

allows fine tuning of alpha for each edge type

# multimodal objective function

$$\mathcal{L}^\alpha(\delta, \mathcal{E}, p) = |\mathcal{E}|^{-1} \left[ \underbrace{\alpha \cdot \sum_{(v, v') \in \mathcal{E}} \delta_{vv'} \cdot p_{vv'}}_{\text{intra-cluster similarity}} + (1 - \alpha) \cdot \underbrace{\sum_{(v, v') \in \mathcal{E}} (1 - \delta_{vv'}) \cdot (1 - p_{vv'})}_{\text{inter-cluster dissimilarity}} \right]$$

would result in giving **more weight** to **dominant type** of edges

$$\mathcal{L}_\beta^\alpha(\delta, \mathcal{E}, p) = \sum_{\substack{x \in \{\mathcal{T}, \mathcal{W}, \mathcal{S}, \mathcal{I}\} \\ y \in \{\mathcal{T}, \mathcal{W}, \mathcal{S}, \mathcal{I}\}}} \beta_{xy} \cdot \mathcal{L}^{\alpha_{xy}}(\delta, \mathcal{E} \cap (x \times y), p) \quad \alpha_{xy} \in [0, 1] \quad \beta_{xy} \in [0, 1] \quad \sum_{x,y} \beta_{xy} = 1$$

**one objective function per edge type subgraph**

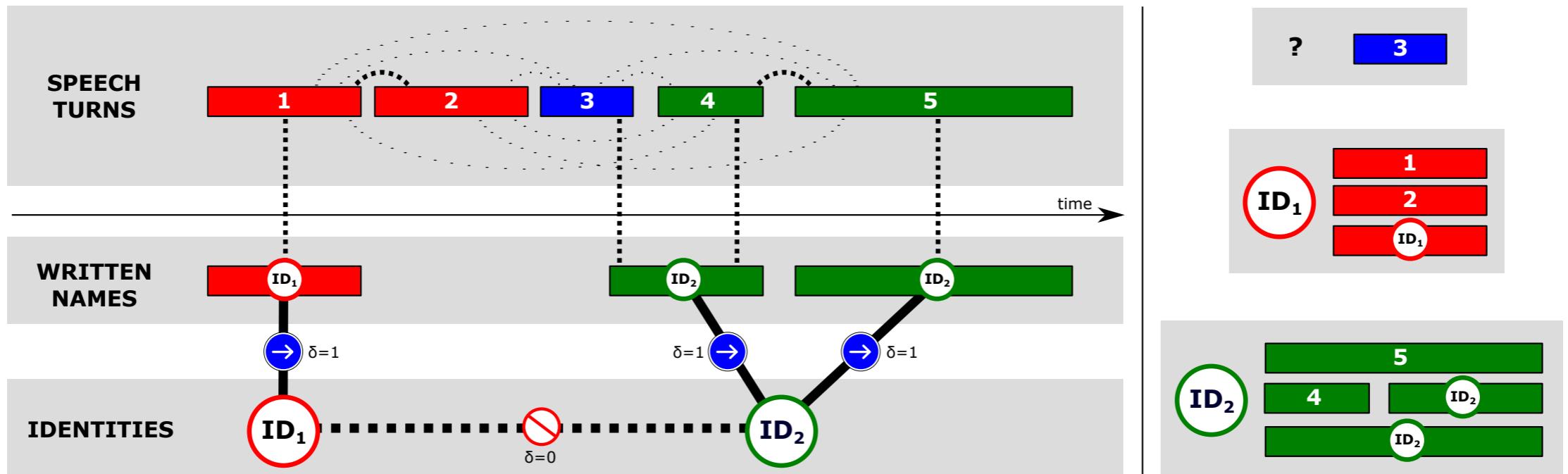
allows different weights beta for different edge types

allows fine tuning of alpha for each edge type

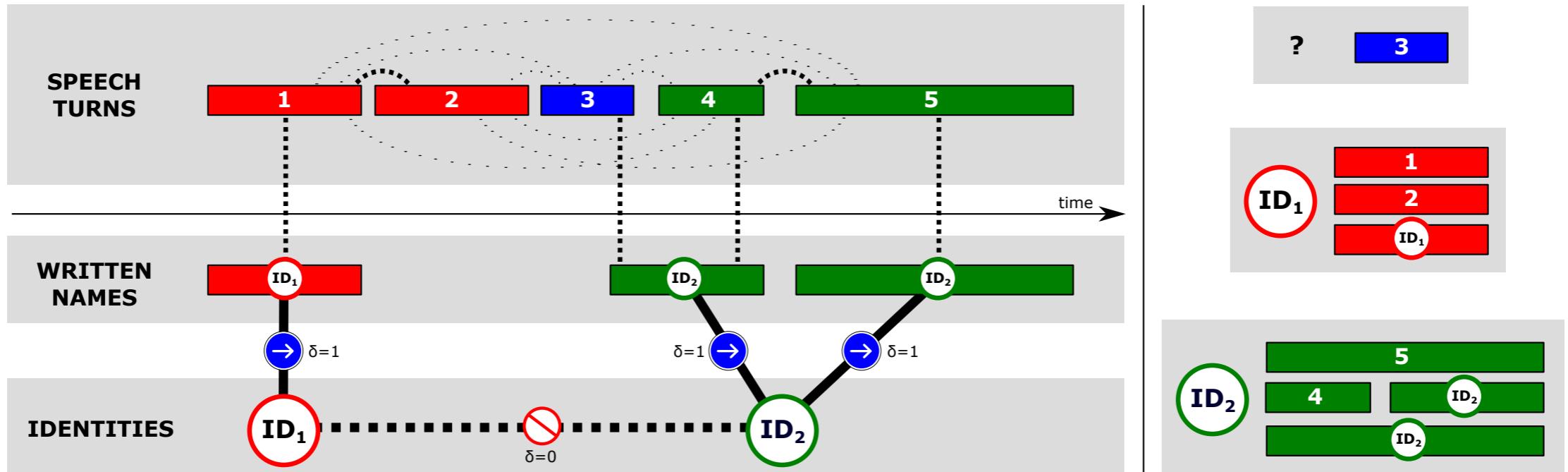
$$\delta^* = \underset{\delta \in \Delta_V}{\operatorname{argmax}} \mathcal{L}_\beta^\alpha(\delta, \mathcal{E}, p)$$

$\alpha$  and  $\beta$  are optimized by random search

# additional constraints



# additional constraints



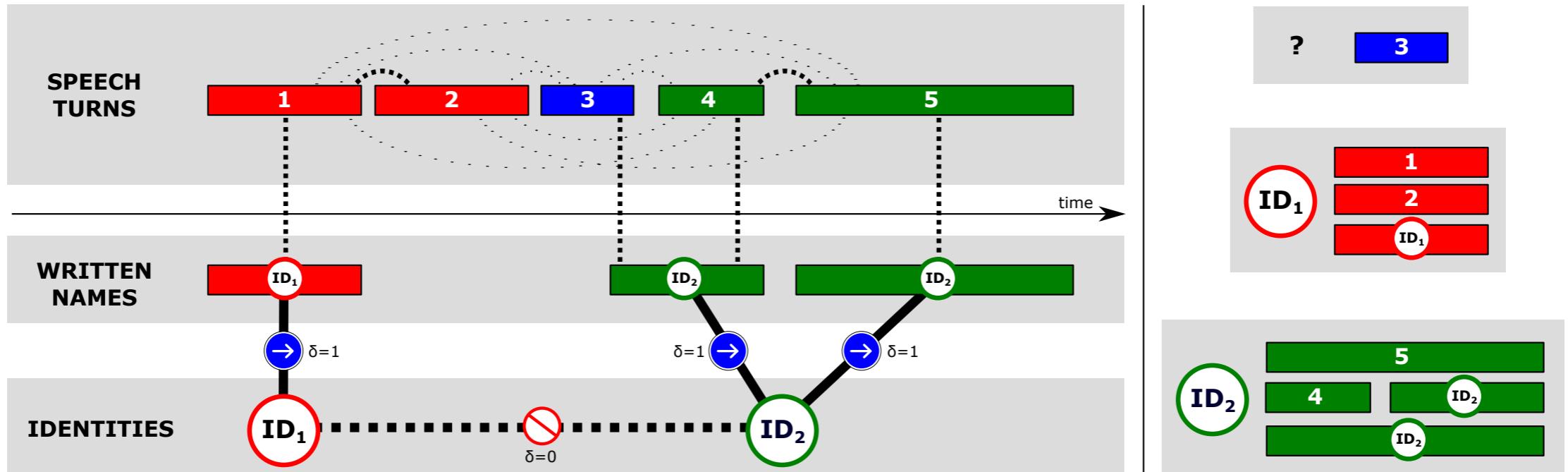
two identity vertices  
cannot end up  
in the same cluster

$$\forall (i, i') \in \mathcal{I}^2, i \neq i' \implies \delta_{ii'} = 0$$

any instance has at most one identity

$$\forall v \in \mathcal{V}, \sum_{i \in \mathcal{I}} \delta_{vi} \leq 1$$

# additional constraints



two identity vertices  
cannot end up  
in the same cluster

$$\forall (i, i') \in \mathcal{I}^2, i \neq i' \implies \delta_{ii'} = 0$$

any instance has at most one identity

$$\forall v \in \mathcal{V}, \sum_{i \in \mathcal{I}} \delta_{vi} \leq 1$$

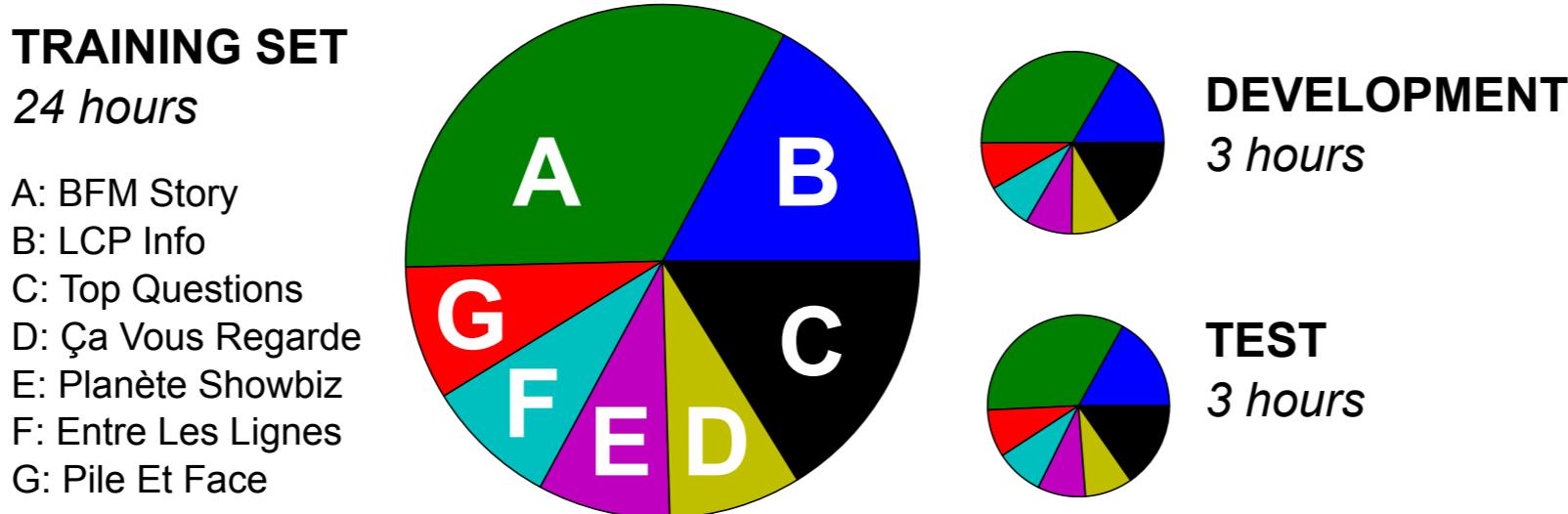
any written or spoken instance  
must be in the same cluster  
as their corresponding identity

$$\forall w \in \mathcal{W}, \delta_{wi_w} = 1$$

$$\forall s \in \mathcal{S}, \delta_{si_s} = 1$$

# experimental protocol

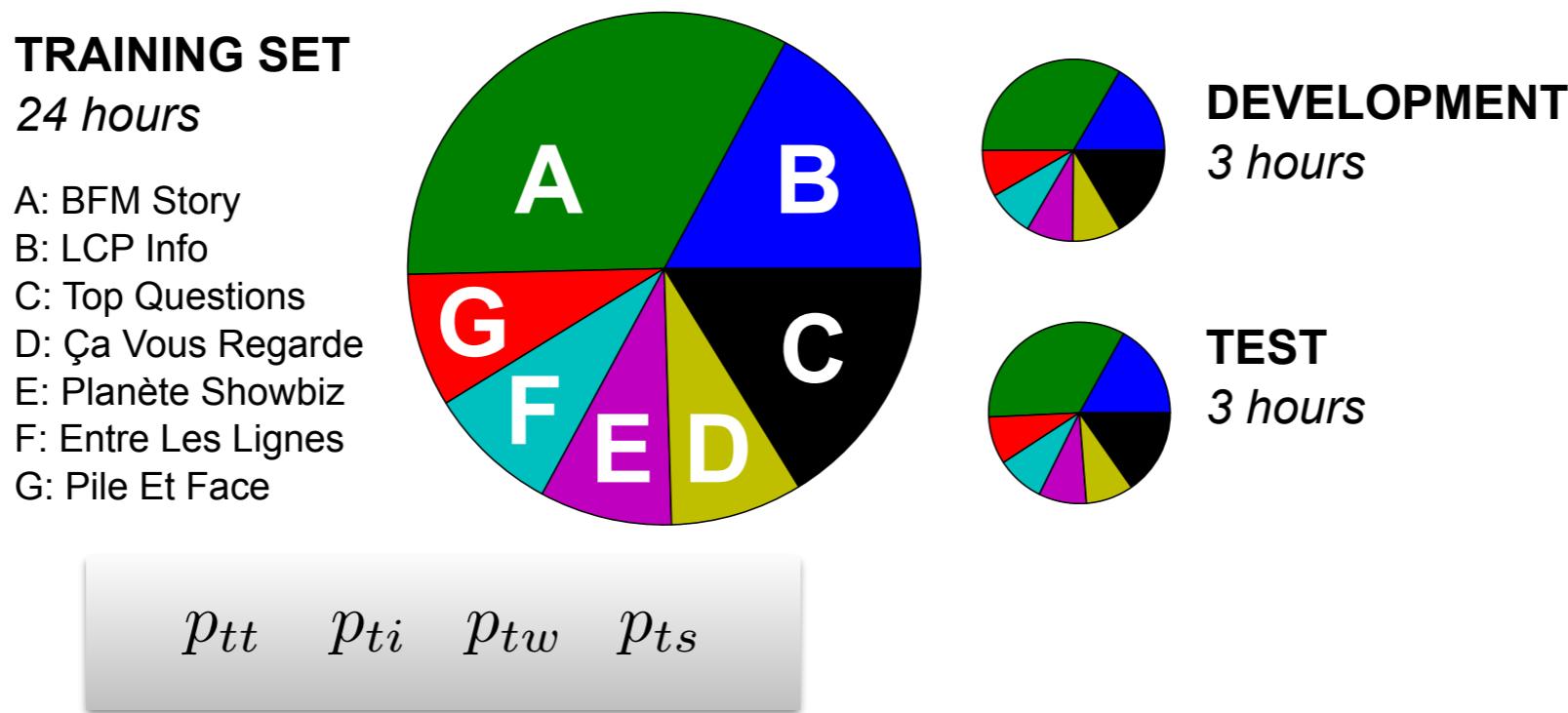
- corpora



- evaluation metrics
  - diarization** error rate
  - identification** error rate / **precision** / **recall**

# experimental protocol

- corpora



- evaluation metrics
  - diarization** error rate
  - identification** error rate / **precision** / **recall**

# experimental protocol

- corpora

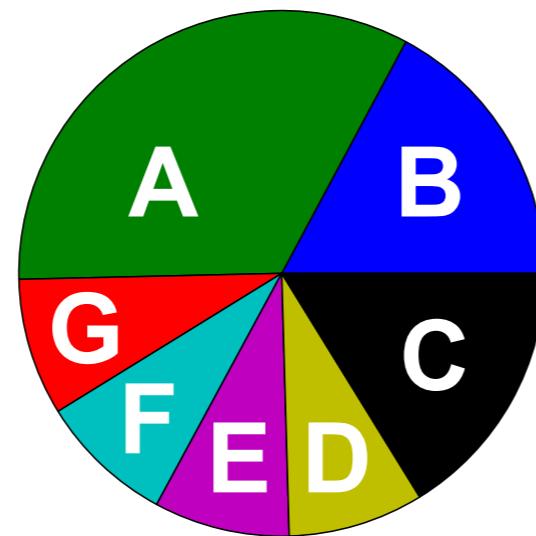
$\text{XER} \in \{\text{IER}, \text{DER}\}$

$$(\alpha^*, \beta^*) = \underset{\alpha, \beta}{\operatorname{argmin}} \mathbb{E}_{\text{dev}} [\text{XER}(r, h)]$$

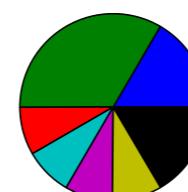
## TRAINING SET

*24 hours*

- A: BFM Story
- B: LCP Info
- C: Top Questions
- D: Ça Vous Regarde
- E: Planète Showbiz
- F: Entre Les Lignes
- G: Pile Et Face

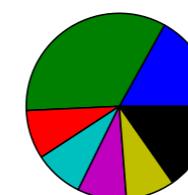


$p_{tt}$     $p_{ti}$     $p_{tw}$     $p_{ts}$



## DEVELOPMENT

*3 hours*



## TEST

*3 hours*

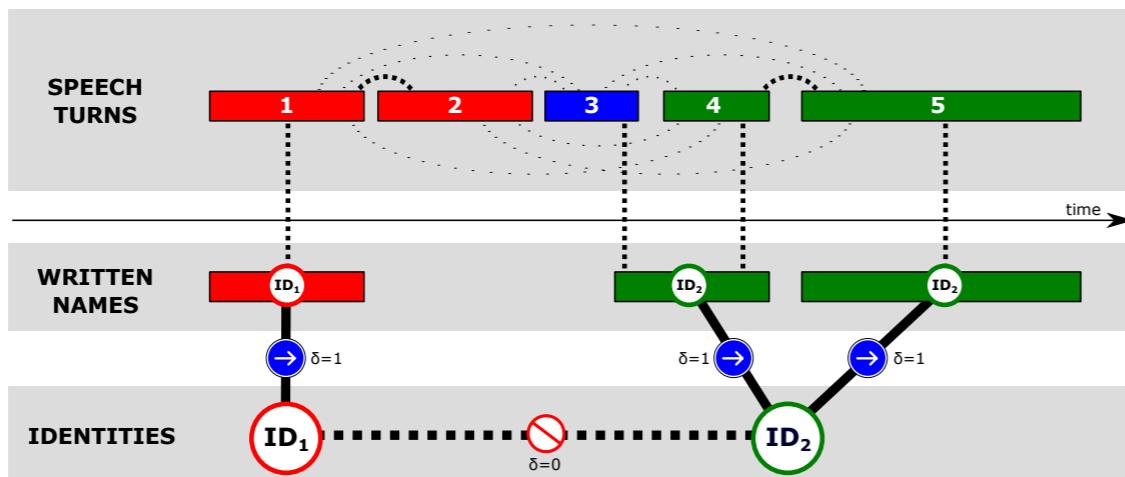
- evaluation metrics
  - diarization** error rate
  - identification** error rate / **precision** / **recall**

# different graph configuration for different **tasks**

- speaker diarization  
speech turn vertices only
- speaker identification
  - **supervised mono-modal** speaker identification  
no written or spoken vertices
  - **unsupervised cross-modal** speaker identification  
no supervised edges  
only written or spoken vertices
  - **multi-modal** speaker identification  
the best of both worlds

# results speaker diarization

Vertices	Edges	DER	Purity	Coverage
$\mathcal{T}$	$t \leftrightarrow t'$	<b>21.1%</b>	<b>94.7%</b>	83.6%
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_W$	$t \leftrightarrow t' \leftrightarrow w \leftrightarrow i_w$	18.3%	93.4%	85.8%
$\mathcal{T} \cup \mathcal{I}^*$	$i^* \leftrightarrow t \leftrightarrow t'$	18.2%	94.0%	86.1%
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_W \cup \mathcal{I}^*$	$i^* \leftrightarrow t \leftrightarrow t' \leftrightarrow w \leftrightarrow i_w$	<b>17.8%</b>	93.9%	85.7%
BIC clustering baseline [1]		<b>19.8%</b>	<b>92.1%</b>	86.8%



written names do help 😊



supervision does help 😊



complementarity 😃



# results mono-modal speaker identification

Vertices	Edges	IER	P.	R.
$\mathcal{T} \cup \mathcal{I}^*$	$i^* \leftrightarrow t$	49.4%	54.7%	54.3%
$\mathcal{T} \cup \mathcal{I}^*$	$i^* \leftrightarrow t \leftrightarrow t'$	<b>47.9%</b>	55.9%	<b>55.8%</b>
GMM-UBM baseline		<b>48.9%</b>	57.5%	54.3%



- the first configuration is equivalent to usual open-set speaker identification

$$\text{ID}(t) = \begin{cases} i^* = \underset{i \in \mathcal{I}^*}{\operatorname{argmax}} p_{ti} & \text{if } p_{ti^*} > (1 - \alpha_{\mathcal{T}\mathcal{I}^*}) \\ \textcircled{?} & \text{otherwise.} \end{cases}$$

- GMM-UBM baseline relies on preliminary speaker diarization step and performs identification at cluster level

- adding  $t \leftrightarrow t'$  edges addresses this limitation

$$\beta_{\mathcal{T}\mathcal{T}} = 0.55 \quad \beta_{\mathcal{T}\mathcal{I}} = 0.45 \quad \alpha_{\mathcal{T}\mathcal{T}} = 0.19$$

# results mono-modal speaker identification

Vertices	Edges	IER	P.	R.
$\mathcal{T} \cup \mathcal{I}^*$	$i^* \leftrightarrow t$	49.4%	54.7%	54.3%
$\mathcal{T} \cup \mathcal{I}^*$	$i^* \leftrightarrow t \leftrightarrow t'$	<b>47.9%</b>	55.9%	<b>55.8%</b>
GMM-UBM baseline		<b>48.9%</b>	57.5%	54.3%



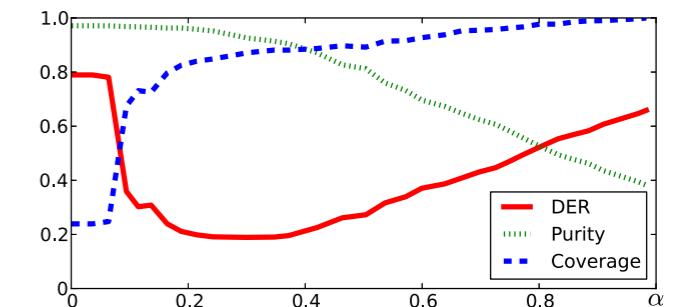
- the first configuration is equivalent to usual open-set speaker identification

$$\text{ID}(t) = \begin{cases} i^* = \underset{i \in \mathcal{I}^*}{\operatorname{argmax}} p_{ti} & \text{if } p_{ti^*} > (1 - \alpha_{\mathcal{T}\mathcal{I}^*}) \\ \text{?} & \text{otherwise.} \end{cases}$$

- GMM-UBM baseline relies on preliminary speaker diarization step and performs identification at cluster level

- adding  $t \leftrightarrow t'$  edges addresses this limitation

$$\beta_{\mathcal{T}\mathcal{T}} = 0.55 \quad \beta_{\mathcal{T}\mathcal{I}} = 0.45 \quad \alpha_{\mathcal{T}\mathcal{T}} = 0.19$$



# results cross-modal speaker identification

Vertices	Edges	IER	Precision	Recall
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_{\mathcal{W}}$	$t \leftrightarrow t' \leftrightarrow w \leftrightarrow i_w$	<b>46.5%</b>	66.8%	56.9%
$\mathcal{T} \cup \mathcal{S} \cup \mathcal{I}_{\mathcal{S}}$	$t \leftrightarrow t' \leftrightarrow s \leftrightarrow i_s$	81.8%	21.5%	21.4%
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_{\mathcal{W}} \cup \mathcal{S} \cup \mathcal{I}_{\mathcal{S}}$	$i_w \Leftrightarrow w \Leftrightarrow t \Leftrightarrow t' \Leftrightarrow s \Leftrightarrow i_s$	<b>45.6%</b>	62.7%	58.2%



# results cross-modal speaker identification

Vertices	Edges	IER	Precision	Recall
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_w$	$t \leftrightarrow t' \leftrightarrow w \leftrightarrow i_w$	<b>46.5%</b>	66.8%	56.9%
$\mathcal{T} \cup \mathcal{S} \cup \mathcal{I}_s$	$t \leftrightarrow t' \leftrightarrow s \leftrightarrow i_s$	81.8%	21.5%	21.4%
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_w \cup \mathcal{S} \cup \mathcal{I}_s$	$i_w \Leftrightarrow w \Leftrightarrow t \Leftrightarrow t' \Leftrightarrow s \Leftrightarrow i_s$	<b>45.6%</b>	62.7%	58.2%

- to be compared with *oracle* performance

# results cross-modal speaker identification

Vertices	Edges	IER	Precision	Recall
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_W$	$t \leftrightarrow t' \leftrightarrow w \leftrightarrow i_w$	<b>46.5%</b>	66.8%	56.9%
$\mathcal{T} \cup \mathcal{S} \cup \mathcal{I}_S$	$t \leftrightarrow t' \leftrightarrow s \leftrightarrow i_s$	81.8%	21.5%	21.4%
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_W \cup \mathcal{S} \cup \mathcal{I}_S$	$i_w \Leftrightarrow w \Leftrightarrow t \Leftrightarrow t' \Leftrightarrow s \Leftrightarrow i_s$	<b>45.6%</b>	62.7%	58.2%

- to be compared with *oracle* performance

Vertices	IER	Precision	Recall
$\mathcal{I}_W$	<b>39.3%</b>	100.0%	<b>63.8%</b>
$\mathcal{I}_S$	35.8%	100.0%	67.3%
$\mathcal{I}_W \cup \mathcal{I}_S$	21.5%	100.0%	81.9%

# results multi-modal speaker identification

Vertices	Edges	All			Anchors			All but anchors		
		IER	P.	R.	IER	P.	R.	IER	P.	R.
$\mathcal{T} \cup \mathcal{I}^*$	$i^* \leftrightarrow t \leftrightarrow t'$	47.9%	55.9%	55.8%	<b>20.3%</b>	<b>86.6%</b>	<b>79.7%</b>	51.8%	50.8%	48.3%
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_{\mathcal{W}}$	$t \leftrightarrow t' \leftrightarrow w \Leftrightarrow i_w$	46.5%	66.8%	56.8%	79.4%	31.8%	20.6%	34.4%	76.9%	65.8%
$\mathcal{T} \cup \mathcal{I}^* \cup \mathcal{W} \cup \mathcal{I}_{\mathcal{W}}$	$i^* \leftrightarrow t \leftrightarrow t' \leftrightarrow w \Leftrightarrow i_w$	<b>25.3%</b>	<b>79.4%</b>	<b>78.6%</b>	23.9%	82.9%	76.1%	<b>22.4%</b>	<b>82.5%</b>	<b>77.9%</b>

# results multi-modal speaker identification

Vertices	Edges	All			Anchors			All but anchors		
		IER	P.	R.	IER	P.	R.	IER	P.	R.
$\mathcal{T} \cup \mathcal{I}^*$	$i^* \leftrightarrow t \leftrightarrow t'$	47.9%	55.9%	55.8%	<b>20.3%</b>	<b>86.6%</b>	<b>79.7%</b>	51.8%	50.8%	48.3%
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_{\mathcal{W}}$	$t \leftrightarrow t' \leftrightarrow w \Leftrightarrow i_w$	46.5%	66.8%	56.8%	79.4%	31.8%	20.6%	34.4%	76.9%	65.8%
$\mathcal{T} \cup \mathcal{I}^* \cup \mathcal{W} \cup \mathcal{I}_{\mathcal{W}}$	$i^* \leftrightarrow t \leftrightarrow t' \leftrightarrow w \Leftrightarrow i_w$	<b>25.3%</b>	<b>79.4%</b>	<b>78.6%</b>	23.9%	82.9%	76.1%	<b>22.4%</b>	<b>82.5%</b>	<b>77.9%</b>

astonished  
face



# results multi-modal speaker identification

Vertices	Edges	All			Anchors			All but anchors		
		IER	P.	R.	IER	P.	R.	IER	P.	R.
$\mathcal{T} \cup \mathcal{I}^*$	$i^* \leftrightarrow t \leftrightarrow t'$	47.9%	55.9%	55.8%	<b>20.3%</b>	<b>86.6%</b>	<b>79.7%</b>	51.8%	50.8%	48.3%
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_W$	$t \leftrightarrow t' \leftrightarrow w \Leftrightarrow i_w$	46.5%	66.8%	56.8%	79.4%	31.8%	20.6%	34.4%	76.9%	65.8%
$\mathcal{T} \cup \mathcal{I}^* \cup \mathcal{W} \cup \mathcal{I}_W$	$i^* \leftrightarrow t \leftrightarrow t' \leftrightarrow w \Leftrightarrow i_w$	<b>25.3%</b>	<b>79.4%</b>	<b>78.6%</b>	23.9%	82.9%	76.1%	<b>22.4%</b>	<b>82.5%</b>	<b>77.9%</b>

astonished  
face



grinning  
face



# results multi-modal speaker identification

Vertices	Edges	All			Anchors			All but anchors		
		IER	P.	R.	IER	P.	R.	IER	P.	R.
$\mathcal{T} \cup \mathcal{I}^*$	$i^* \leftrightarrow t \leftrightarrow t'$	47.9%	55.9%	55.8%	<b>20.3%</b>	<b>86.6%</b>	<b>79.7%</b>	51.8%	50.8%	48.3%
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_W$	$t \leftrightarrow t' \leftrightarrow w \Leftrightarrow i_w$	46.5%	66.8%	56.8%	79.4%	31.8%	20.6%	34.4%	76.9%	65.8%
$\mathcal{T} \cup \mathcal{I}^* \cup \mathcal{W} \cup \mathcal{I}_W$	$i^* \leftrightarrow t \leftrightarrow t' \leftrightarrow w \Leftrightarrow i_w$	<b>25.3%</b>	<b>79.4%</b>	<b>78.6%</b>	23.9%	82.9%	76.1%	<b>22.4%</b>	<b>82.5%</b>	<b>77.9%</b>

astonished  
face



grinning  
face



- **the best of both worlds**

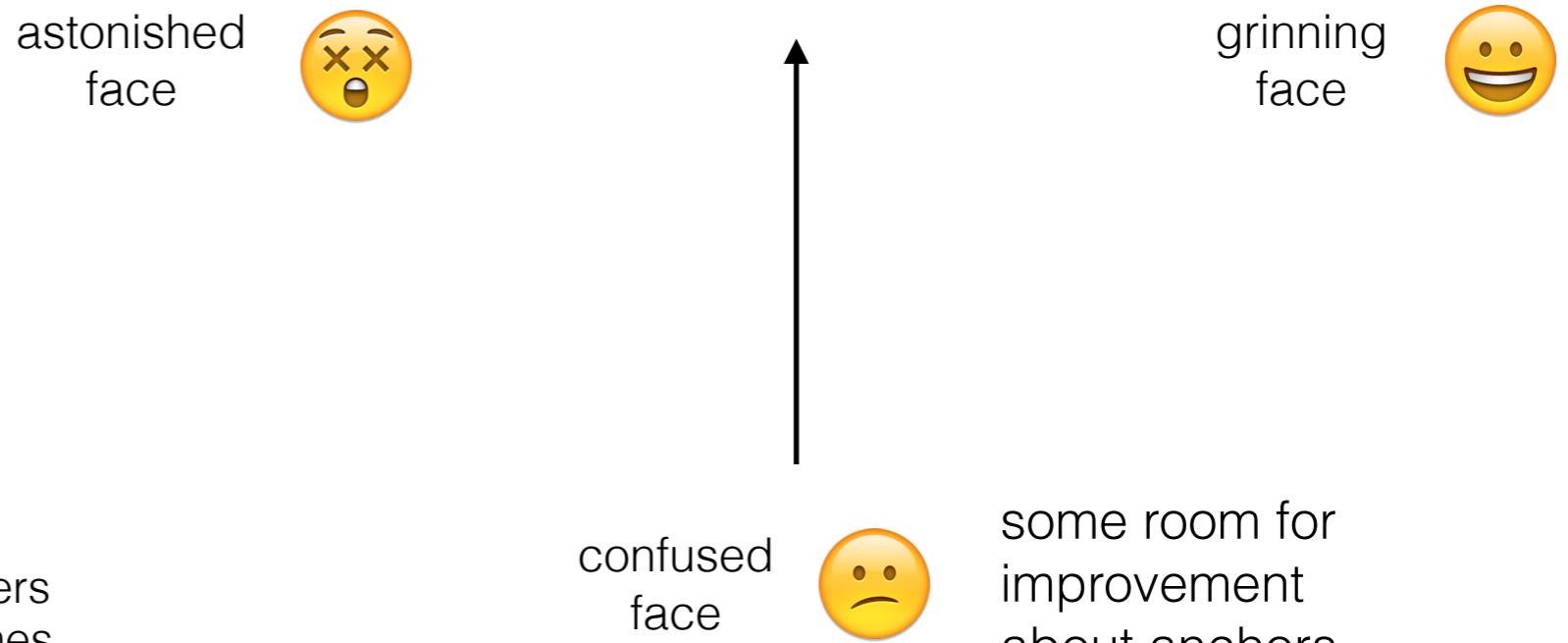
- supervised for anchors  
lots of training data
- cross-modal for guests or reporters  
often introduced by overlaid names
- **com-ple-men-ta-ry** approaches

# results multi-modal speaker identification

Vertices	Edges	All			Anchors			All but anchors		
		IER	P.	R.	IER	P.	R.	IER	P.	R.
$\mathcal{T} \cup \mathcal{I}^*$	$i^* \leftrightarrow t \leftrightarrow t'$	47.9%	55.9%	55.8%	<b>20.3%</b>	<b>86.6%</b>	<b>79.7%</b>	51.8%	50.8%	48.3%
$\mathcal{T} \cup \mathcal{W} \cup \mathcal{I}_W$	$t \leftrightarrow t' \leftrightarrow w \Leftrightarrow i_w$	46.5%	66.8%	56.8%	79.4%	31.8%	20.6%	34.4%	76.9%	65.8%
$\mathcal{T} \cup \mathcal{I}^* \cup \mathcal{W} \cup \mathcal{I}_W$	$i^* \leftrightarrow t \leftrightarrow t' \leftrightarrow w \Leftrightarrow i_w$	<b>25.3%</b>	<b>79.4%</b>	<b>78.6%</b>	23.9%	82.9%	76.1%	<b>22.4%</b>	<b>82.5%</b>	<b>77.9%</b>

- **the best of both worlds**

- supervised for anchors  
lots of training data
- cross-modal for guests or reporters  
often introduced by overlaid names
- **com-ple-men-ta-ry** approaches



some room for improvement about anchors

# dem



Camomile @ ERRARE

localhost:8070/#/diff

Connected as bredin Logout

Camomile Difference Regression Fusion Segmentation

Corpus: REPERE/phase2/test Medium: BFMTV\_BFMStory\_2012-07-24\_175800

Reference: SPEAKER (on 2014-02-14) Hypothesis: QCOMPERE Primary Supervised (speaker)



RONALD GUINTRANGE  
ENVOYÉ SPÉCIAL À LONDRES (GRANDE-BRETAGNE)  
DIRECT 18:18  
LONDRES J-3  
FRANCE Expertise de Bercy sur la vente de l'hippodrome de CAC 40 0:20:57

00:16:00 00:17:00 00:18:00 00:19:00 00:20:00 00:21:00 00:22:00 00:23:00

Reference

Hypothesis

Difference

00:15:00 00:30:00 00:45:00

Display piechart

Camomile Project

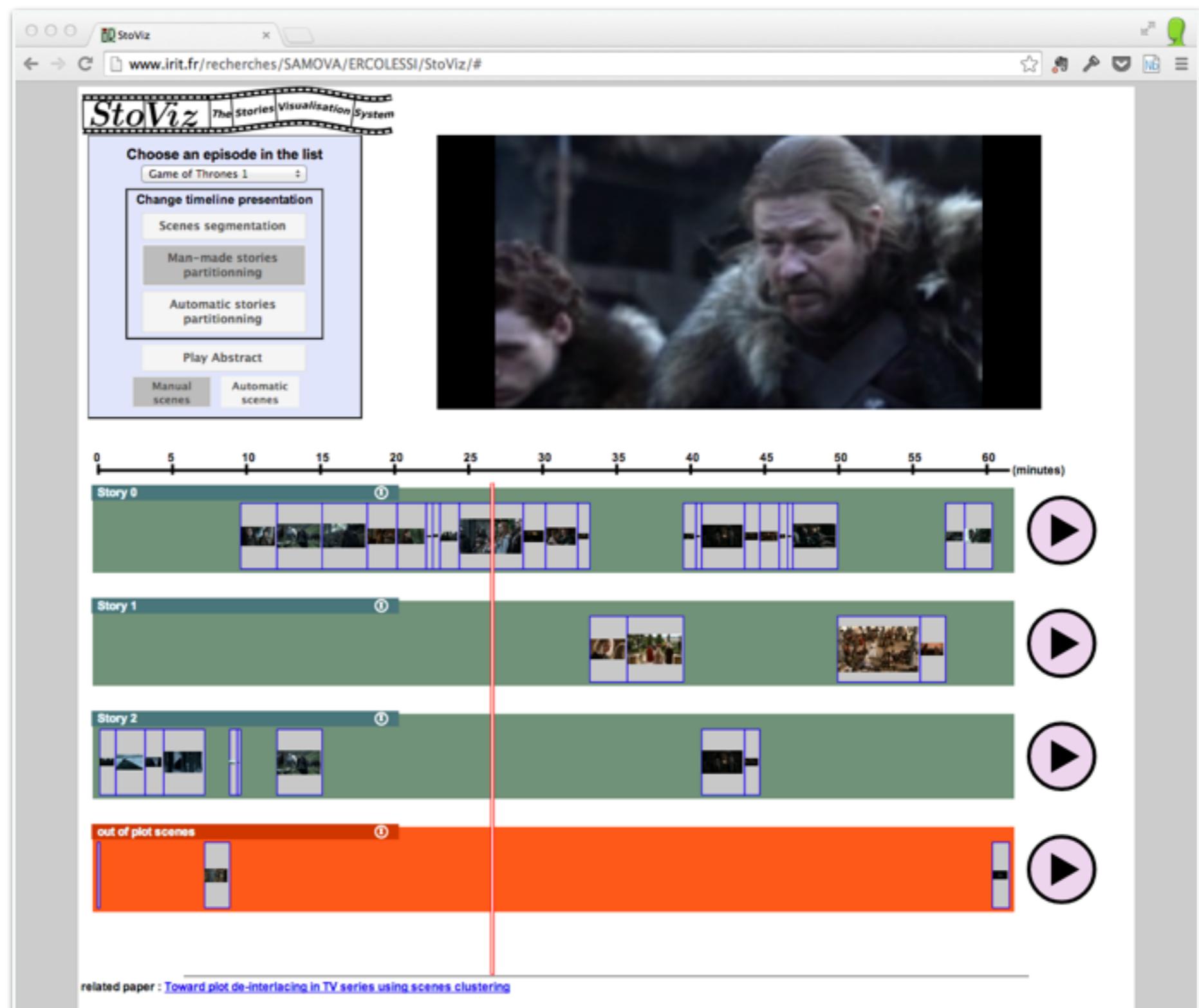
# conclusion

- a global **optimisation framework** for **person identification** in multimedia documents
- a **unified graphical** approach for mono-, cross- and multi-modal identification **easily extendable** to additional modalities
- **main limitations — scalability** problem quickly becomes intractable as the number of vertices increases

# what's next?

- extension to **face** modality
  - $f \leftrightarrow f'$  face clustering
  - $f \leftrightarrow i$  supervised face recognition
  - $t \leftrightarrow f$  audio-based face recognition
- **alternative** objective functions  
integrate **prior knowledge** about the temporal **structure** of interactions
- **named speaker identification**  
better use of spoken names when they are the only source of information
- other **multimodal clustering** problems  
scene clustering for TV series plot de-interlacing

# stoviz





# person instance **graphs** for person **recognition** in **multimedia** data

*Hervé Bredin, Anindya Roy, Viet-Bac Le and Claude Barras*

**Person Instance Graphs for Mono-, Cross- and Multi-Modal Person Recognition in Multimedia Data:  
Application to Speaker Identification in TV Broadcast**

International Journal of Multimedia Information Retrieval, May 2014

*Hervé Bredin, Antoine Laurent, Achintya Sarkar, Viet-Bac Le, Claude Barras and Sophie Rosset*

**Person Instance Graphs for Named Speaker Identification**

Odyssey 2014: the Speaker and Language Recognition Workshop, June 2014



*Anindya Roy, Camille Guinaudeau, Hervé Bredin and Claude Barras*

**TVD: a Reproducible and Multiply Aligned TV Series Dataset**

LREC 2014: the 9th edition of the Language Resources and Evaluation Conference, May 2014

**tvd.niderb.fr**