

# Automatique Evaluation of Natural Language Generation using Measures of Similarity

Pierre Colombo

LEYA Lab 16/12/2021

# What is automatic evaluation?

**Problem:**

$S_1$ : The weather is cold today.

$S_2$ : It is freezing today



0.8

**Similar**

$S_1$ : I like those cats.

$S_2$ : It is freezing today



0.1

**Dissimilar**

**Goal:**

Building a metric  $m$

$$m : \mathcal{S} \times \mathcal{S} \rightarrow [0,1]$$

$$(S_1, S_2) \rightarrow m(S_1, S_2)$$

**Success Criterion:**

When do we know that  $m$  is good?



**Correlation with human scores**

Koehn 2009; Specia, Raj, and Turchi 2010; Chatzikoumi 2020

# Importance of Evaluation of NLG

---

## Why is automatic evaluation popular?

1. **Cheap**: compared to human evaluation.
2. **Fast**: you can label “instantaneously”.
3. **Reproducible**: two sentences always get the same score.
4. **Easy** to use (e.g no annotator training, no form design).

Karpinska et al. 2021

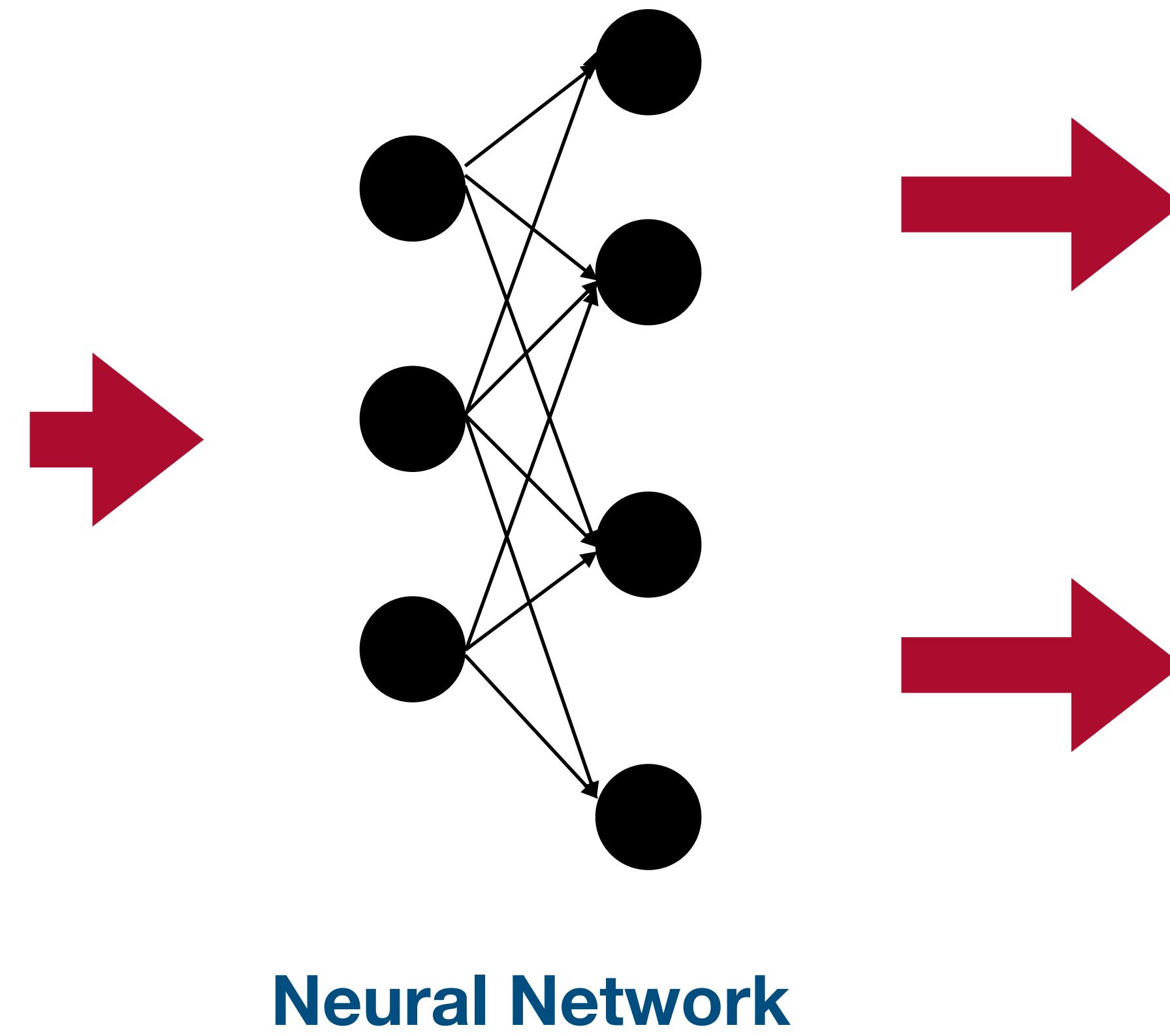
## Why do we need evaluation of NLG?

1. **Debug** NLG systems without annotators.
2. **Improve** learning of systems by deriving new losses.
3. **Compare** different systems.

# Statistical Measures of Similarity

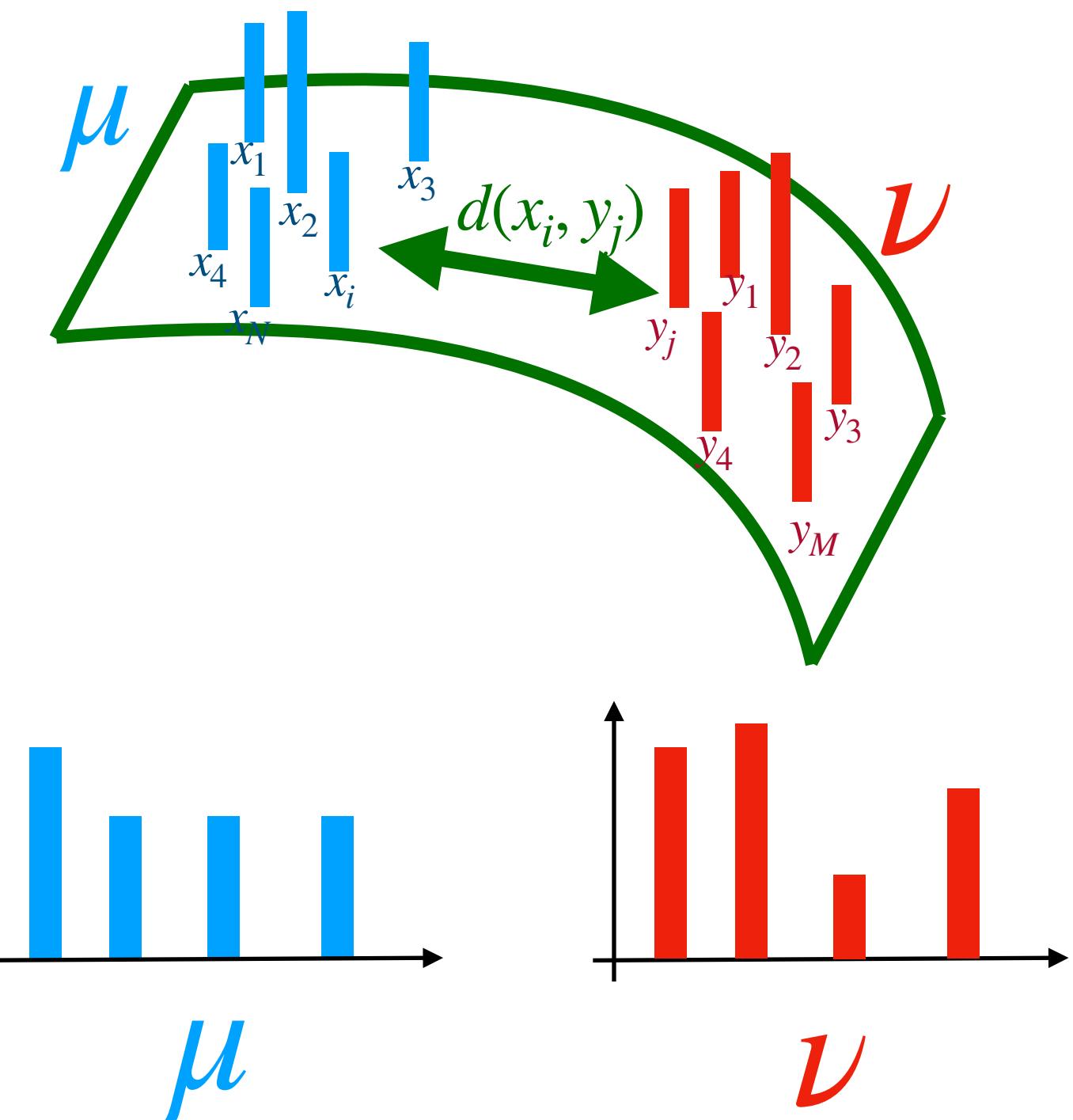
Hello, Chicago.  
If there is anyone out  
there who still doubts  
that America is a place  
where all things are  
possible, who still  
wonders if the dream of  
our founders is alive in  
our time, [...].  
Yes we can!

Input Text



Neural Network

## High dimensional data



Soft Probabilities

When working with Neural Networks, we need to compare probability distributions

The statistical measures are a tool measure this similarity/dissimilarity between probability distribution!

# Outline

---

## 1. Context: examples of problems, evaluation of automatic evaluation

### 1.1 Types of Automatic Evaluation of NLG

### 1.2 Evaluation of automatic evaluation

## 2. Reference Based Automatic Evaluation

### 2.1 DepthScore

Guillaume Staerman, Pavlo Mozharovskyi, Pierre Colombo, Stephan Cléménçon, Florence d'Alché-Buc. A Pseudo-Metric between Probability Distributions based on Depth-Trimmed Regions.

### 2.2 BaryScore

Pierre Colombo, Guillaume Staerman, Chloé Clavel, Pablo Piantanida. Automatic Text Evaluation through the Lens of Wasserstein Barycenters. EMNLP 2022 (oral)

### 2.2 InfoLM

Pierre Colombo, Chloé Clavel and Pablo Piantanida. InfoLM: A New Metric to Evaluate Summarization & Data2Text Generation. AAAI 2022

## 3. Conclusions

# **1. Context: examples of problems, evaluation of automatic evaluation.**

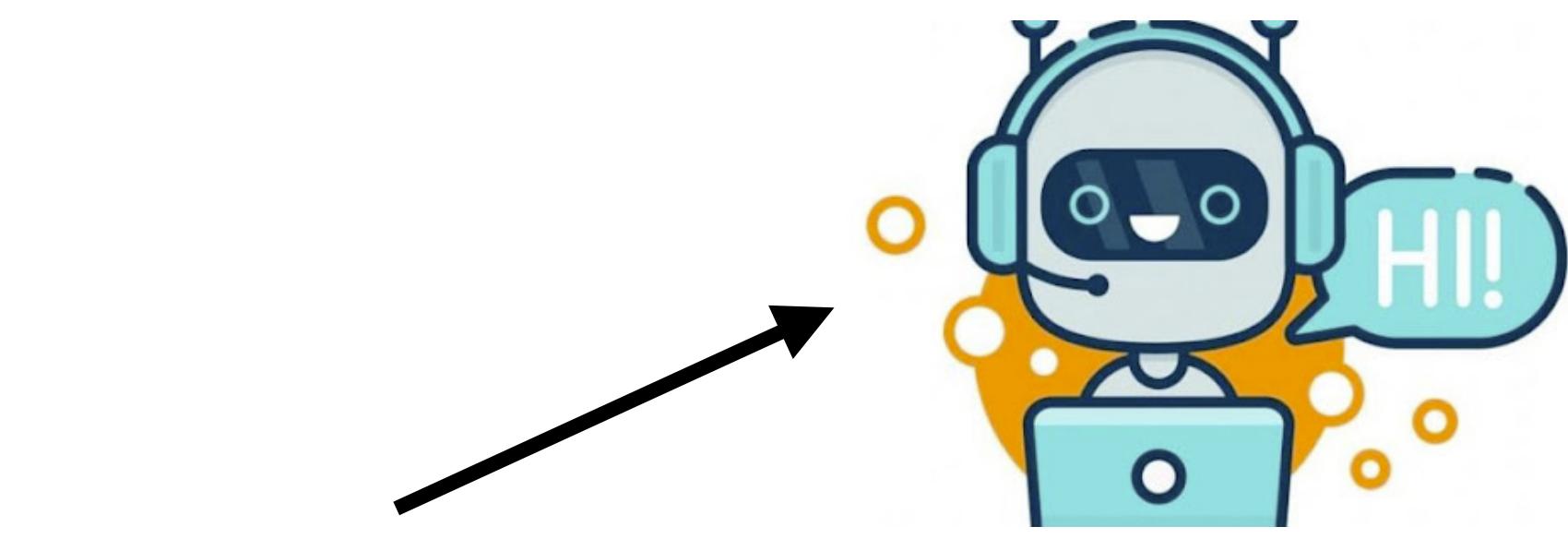
## Reference based vs reference free evaluation

$$m : \mathcal{S} \times \mathcal{S} \rightarrow [0,1]$$

$$(S_1, S_2) \rightarrow m(S_1, S_2)$$

What are  $S_1$  and  $S_2$  ?

Reference based



I like running outside !

System



Human

$S_1$

J'aime courir à l'extérieur !

$S_2$

J'apprécie courir dehors !

Reference free

$S_1$  I like running outside !



$S_2$  J'aime courir à l'extérieur !

# Sentence level vs word level scores

## Sentence Level Scores

**Goal:**

**Building a metric  $m$**

$$m : \mathcal{S} \times \mathcal{S} \rightarrow [0,1]$$

$$(S_1, S_2) \rightarrow m(S_1, S_2)$$

## Word Level Scores: Word-level quality estimation

Ding et al. 2021

Rei et al. 2021

Target	Vaccination	put	end	to	the	pandemic	Space
Post-edit	Vaccination	put	an	ends	to	the	pandemic
TER	OK	OK	OK	BAD	BAD	OK	OK OK

# Metric Evaluation

**Success Criterion:** Correlation with human score

Notations	S systems N texts	$R_i$ i-th reference	$h(C_i^j)$ human score
		$C_i^j$ i-th text candidate generated by j-th system	

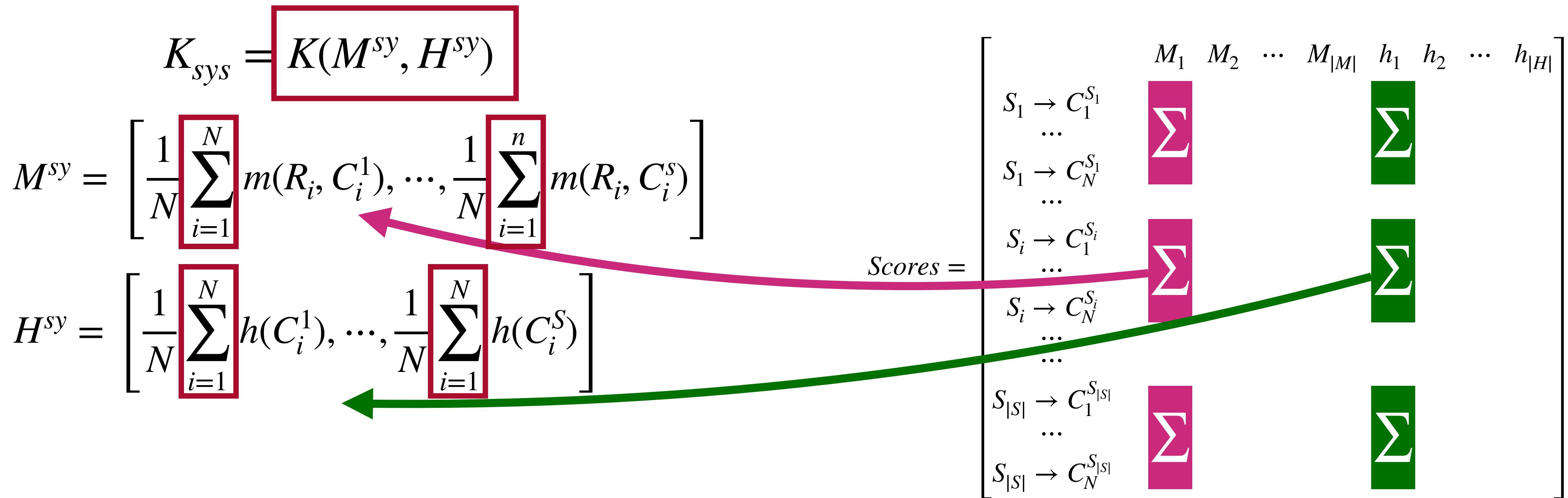
$$Scores = \begin{bmatrix} & M_1 & M_2 & \cdots & M_{|M|} & h_1 & h_2 & \cdots & h_{|H|} \\ S_1 \rightarrow C_1^{S_1} & & & & & & & & \\ \cdots & & & & & & & & \\ S_1 \rightarrow C_N^{S_1} & & & & & & & & \\ \cdots & & & & & & & & \\ \cdots & & M_{i_0}(C_j^i, R_j) & & & & h_{k_0}(C_j^i, R_j) & & \\ S_{|S|} \rightarrow C_1^{S_{|S|}} & & & & & & & & \\ \cdots & & & & & & & & \\ S_{|S|} \rightarrow C_N^{S_{|S|}} & & & & & & & & \end{bmatrix}$$

Can the metric be used to compare  
the performance of two systems?

Can the metric be used as a loss or reward  
of a system?

# Metric Evaluation: system level aggregation

Can the metric be used to compare the performance of two systems?



System Aggregation !  
Compare vector of length  $S$

# Metric Evaluation: text level Aggregation

Can the metric be used as a loss or reward of a system?

$$K_{text} = \frac{1}{N} \sum_{i=1}^N K(M_i^{text}, H_i^{text})$$

$$H_i^{text} = [h(C_i^1), \dots, h(C_i^S)]$$

$$M_i^{text} = [m(R_i, C_i^1), \dots, m(R_i, C_i^S)]$$

$$Scores = \begin{bmatrix} & M_1 & M_2 & \cdots & M_{|M|} & h_1 & h_2 & \cdots & h_{|H|} \\ S_1 \rightarrow C_1^{S_1} & \text{[Red]} & & & & \text{[Green]} & & & \text{[Green]} \\ \dots & & & & & & & & \\ S_1 \rightarrow C_N^{S_1} & \text{[Red]} & & & & & & & \text{[Green]} \\ \dots & & & & & & & & \\ S_i \rightarrow C_1^{S_i} & \text{[Red]} & & & & & & & \text{[Green]} \\ \dots & & & & & & & & \\ S_i \rightarrow C_N^{S_i} & \text{[Red]} & & & & & & & \\ \dots & & & & & & & & \\ S_{|S|} \rightarrow C_1^{S_{|S|}} & \text{[Red]} & & & & & & & \text{[Green]} \\ \dots & & & & & & & & \\ S_{|S|} \rightarrow C_N^{S_{|S|}} & \text{[Red]} & & & & & & & \text{[Green]} \end{bmatrix}$$

Text Aggregation !  
Averaged correlation

## **A summary so far .....**

SUMMARY

### **What are the different settings?**

**Reference Based/ Reference Free**

**Different Granularities (word vs sentence)**



### **How to evaluate a metric?**

**Text Level Correlation**

**Can the metric be used as a loss or reward of a system?**

**System Level Correlation**

**Can the metric be used to compare the performance of two systems?**

**Let's speak about the different metrics for reference-based evaluation!**

## 2. Reference Based Automatic Evaluation of NLG

Pierre Colombo, Guillaume Staerman, Chloé Clavel, Pablo Piantanida. Automatic Text Evaluation through the Lens of Wasserstein Barycenters. (oral) EMNLP 2021

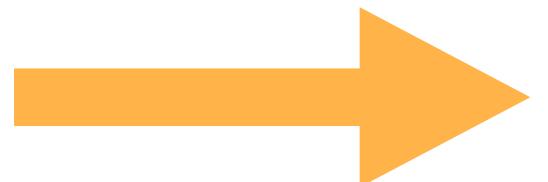
Pierre Colombo, Chloé Clavel and Pablo Piantanida. InfoLM: A New Metric to Evaluate Summarization & Data2Text Generation. AAAI 2022

Guillaume Staerman, Pavlo Mozharovskyi, Pierre Colombo, Stephan Cléménçon, Florence d'Alché-Buc. A Pseudo-Metric between Probability Distributions based on Depth-Trimmed Regions. Submitted

# Existing Metrics for Reference Based NLG

**Goal**

R: The weather is cold today.  
C: It is freezing today



0.8

**Similar**

R: I like those cats.  
C: It is freezing today



0.1

**Dissimilar**

**Edit Based**

**N-gram Based**

**Embedding Based**

# Existing Methods

## Edit Based

Snover et al. 2006

### Operations

- Insertion (I)
- Deletion (D)
- Substitution (S).

tailor -> sailor (**S**)

sailor -> sailir (**S**)

sailir -> sailin (**S**)

sailin\_ -> sailing (**I**)

Distance is 4 !

## N-gram Based

Papineni et al. 2002

C : I like these very nice pies !

R : I like those cakes !

### Unigrams

C : I like these very nice pies !

R : I like those cakes !

### Bigrams

C : I like these very nice pies !

R : I like those cakes !

## Embedding Based

### Word Mover distance

Kusner et al. 2015

### BertScore

Zhang et al. 2019

### MoverScore

Zhao et al. 2019

### Sentence Mover

Clark et al. 2019

## Existing Methods

### Edit Based

Snover et al. 2006

#### Operations

- Insertion (I)
- Deletion (D)
- Substitution (S).

tailor -> sailor (**S**)

sailor -> sailir (**S**)

sailir -> sailin (**S**)

sailin\_ -> sailing (**I**)

Distance is 4 !

### N-gram Based

Papineni et al. 2002

C : I like these very nice pies !

R : I like those cakes !

#### Unigrams

C : I like these very nice pies !

R : I like those cakes !

#### Bigrams

C : I like these very nice pies !

R : I like those cakes !

## BertScore & DepthScore

### Embedding Based

#### Word Mover distance

Kusner et al. 2015

#### BertScore

Zhang et al. 2019

#### MoverScore

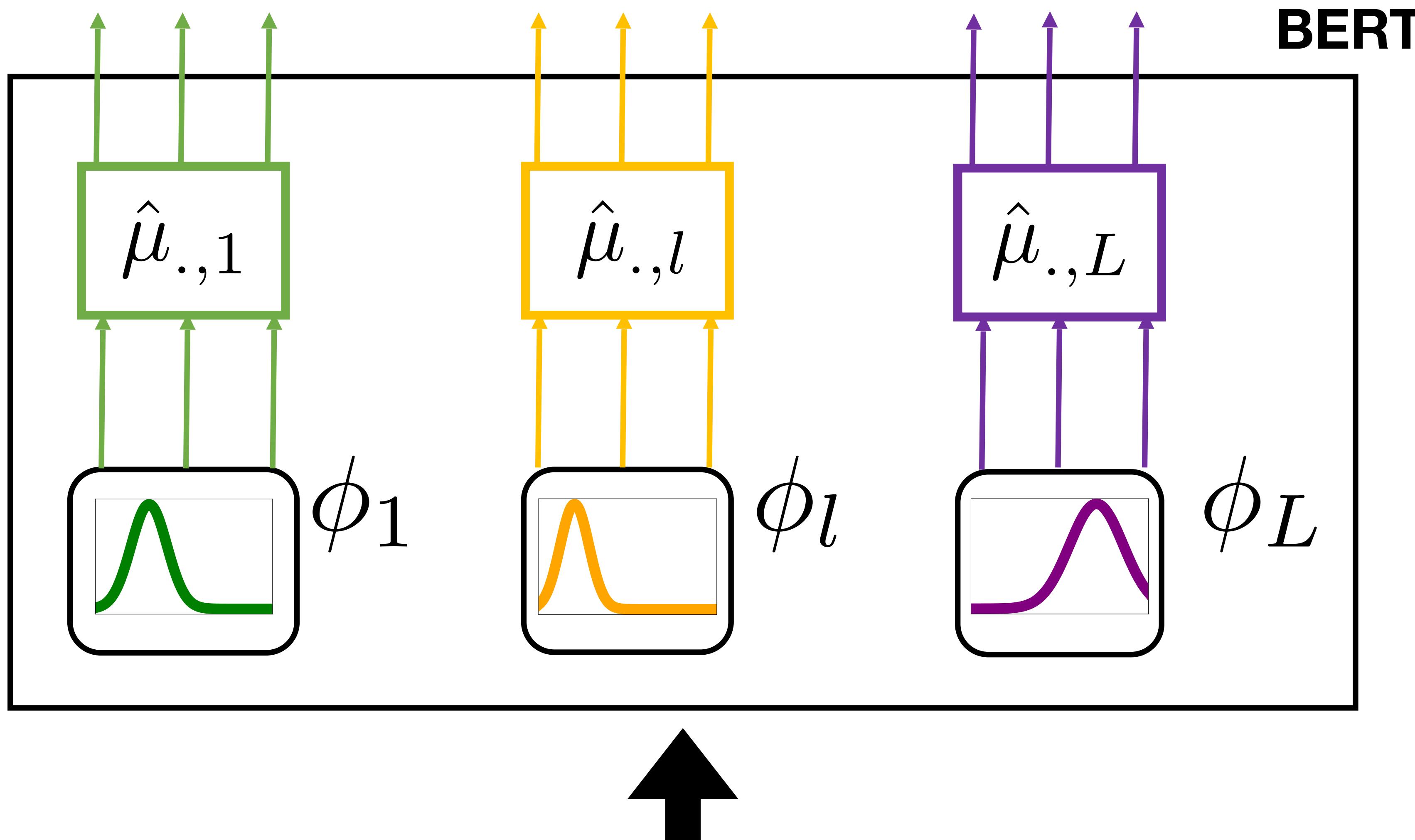
Zhao et al. 2019

#### Sentence Mover

Clark et al. 2019

# A recall on the hidden space of BERT

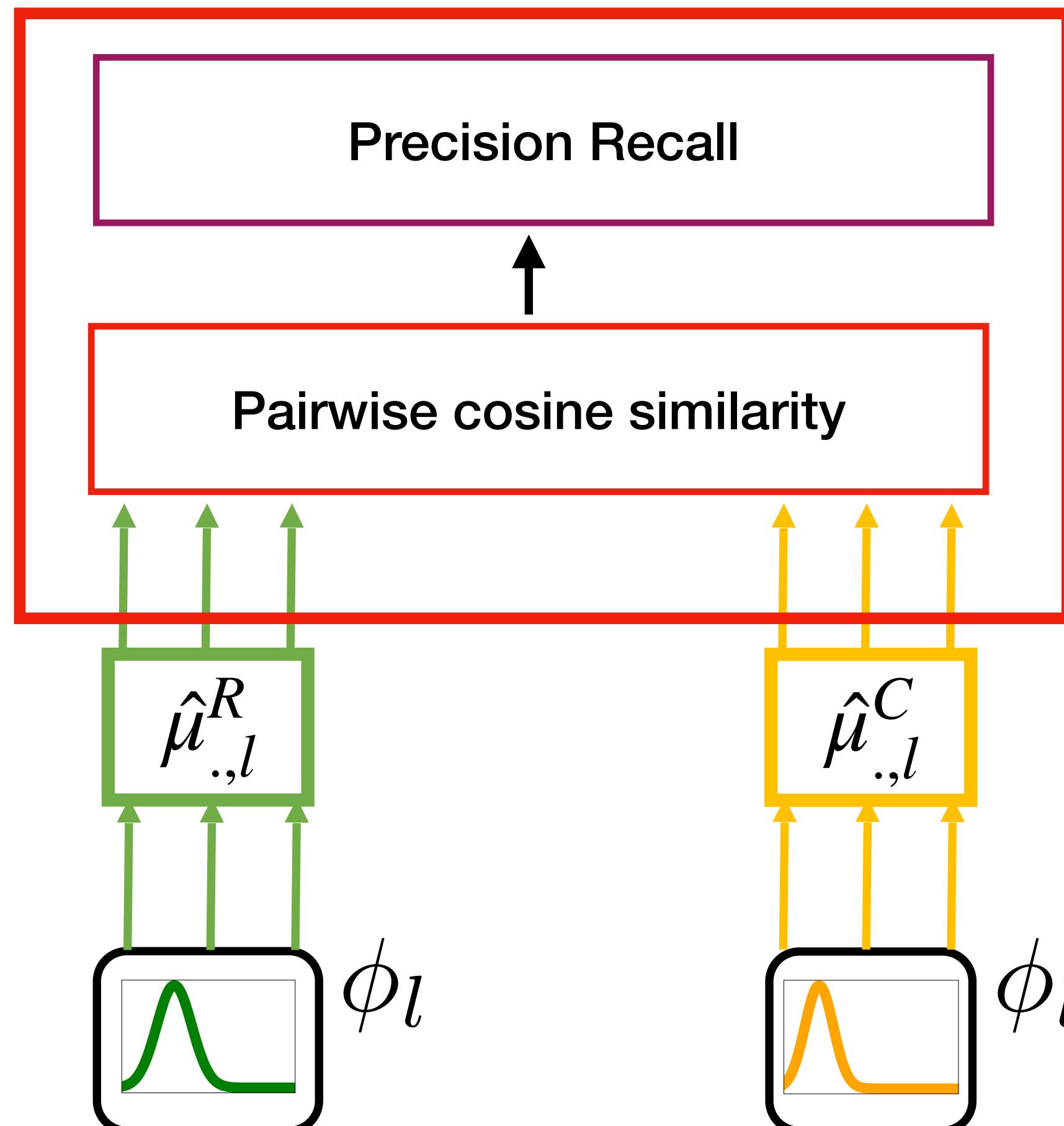
Delvin et al. 2019



1. “Contextual Embedding”
2. 12 Layers ( $L=12$ )
3.  $H = 768$
4. Max length = 512

# BertScore

Zhang et al. 2019



This is a distance between two empirical distributions !

## Advantage

1. Deal with **paraphrases**
2. Include “**semantic**”

## Limitations

1. Use only **one layer**
2. Use **arbitrary sequence of operation**

R: The weather  
is cold today

C: It is freezing  
this morning

# DepthScore

Guillaume Staerman, Pavlo Mozharovskyi, Pierre Colombo, Stephan Clémençon, Florence d'Alché-Buc.  
A Pseudo-Metric between Probability Distributions based on Depth-Trimmed Regions. Submitted

# A glance on data-depths

G. Staerman, P. Mozharovskyi, P. Colombo, S. Cléménçon, F. d'Alché-Buc. A Pseudo-Metric between Probability Distributions based on Depth-Trimmed Regions.

Let's build an alternative distance between distribution !

## Definition

Data depth are non parametric statistic which measures the **centrality** of any element of a  $\mathcal{X}$  w.r.t. a **probability distribution**  $\mathcal{P}_{\mathcal{X}}$ .

$$D : \mathcal{X} \times \mathcal{P}(\mathcal{X}) \rightarrow [0,1]$$

$$(x, \rho) \rightarrow D(x, \rho) = D_{\rho}(x)$$

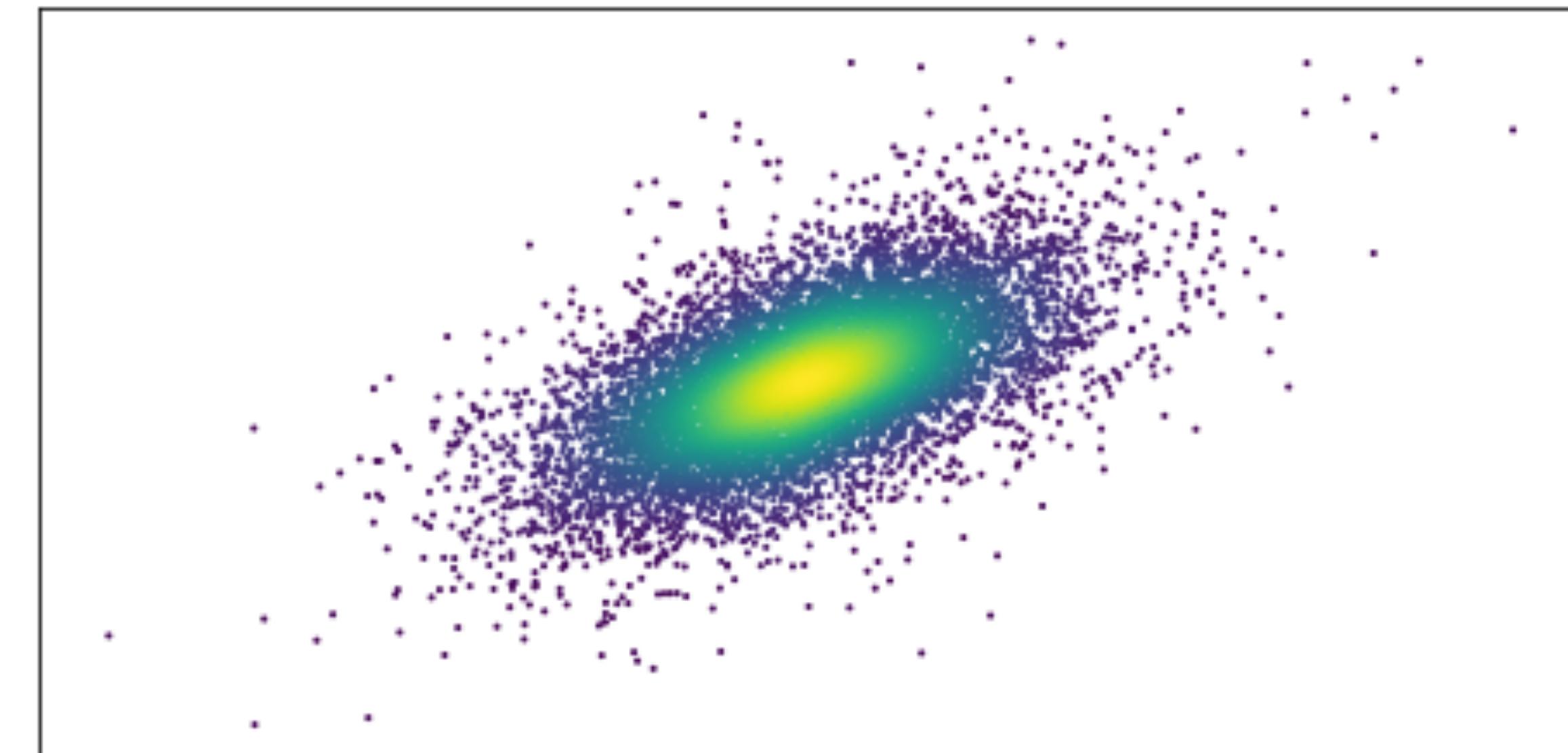
## Properties

Affine Invariant

Monotonous on rays

Maximal at symmetry center

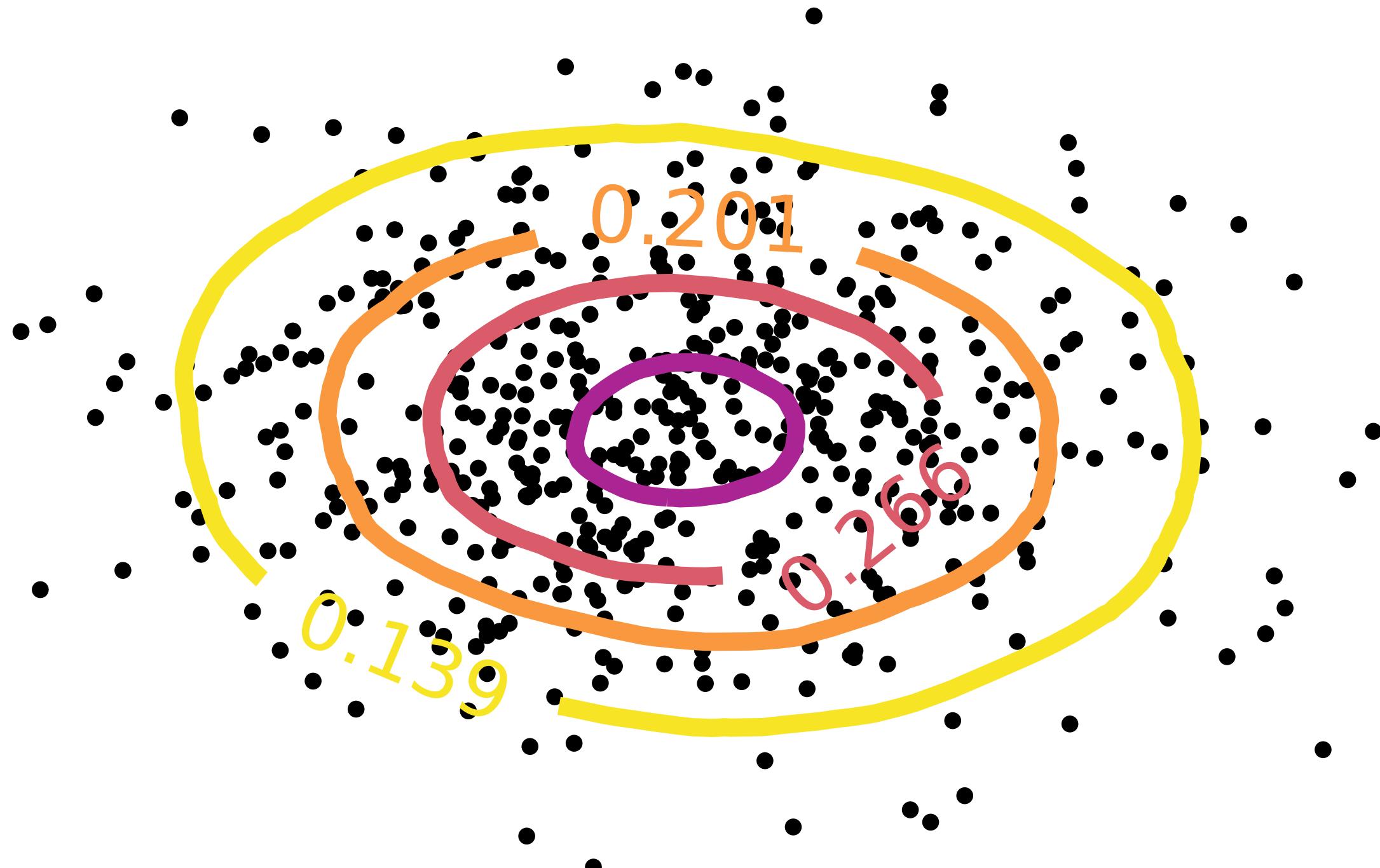
Vanishing at infinity



# From data-depths to a distance

**Definition** For any  $\alpha \in [0,1]$ , the associated  $\alpha$ -depth region of a depth function is defined as:

$$D_\rho^\alpha = \left\{ x \in \mathbb{R}^d, D_\rho(x) \geq \alpha \right\}$$



## Notes

Depth regions are nested and they generalise the notion of quantiles to a multivariate distribution.

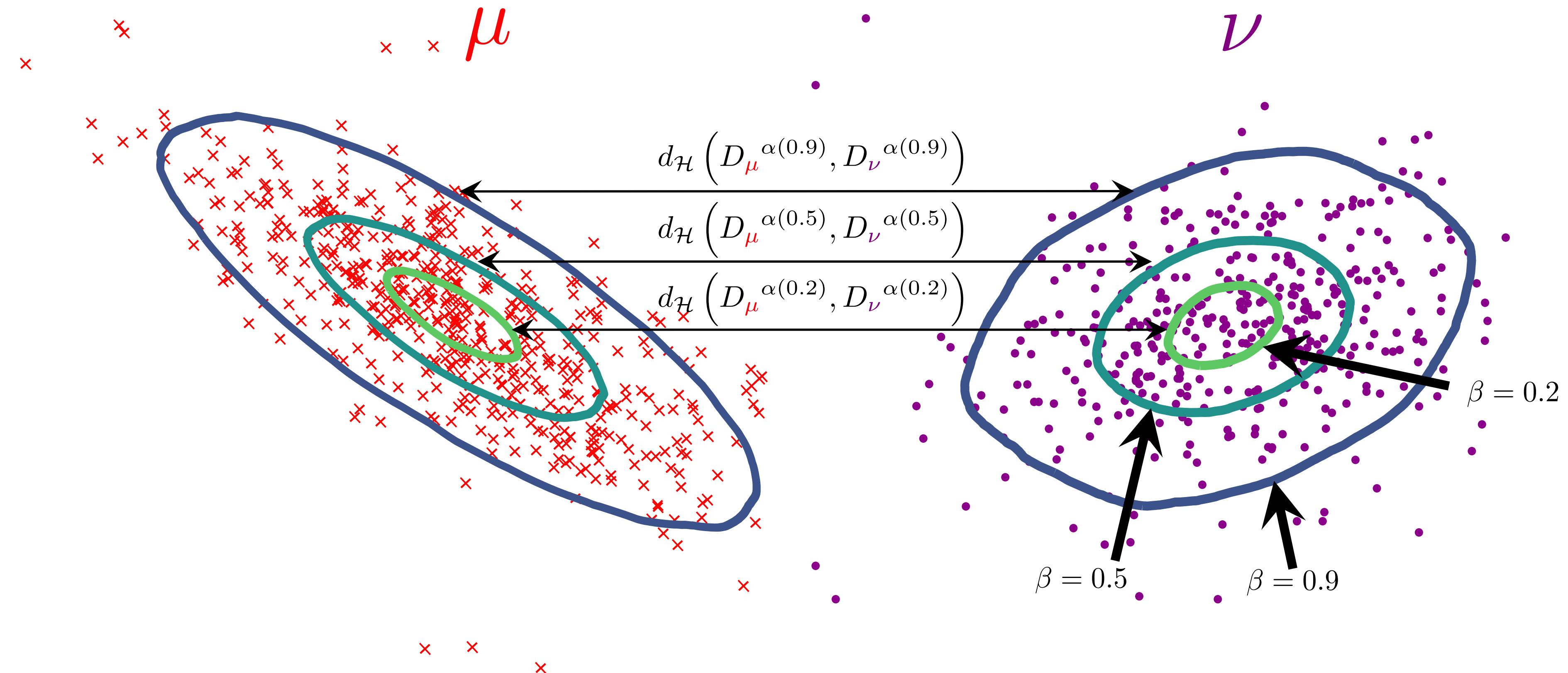
# A pseudo metric based on depth-trimmed regions

**Definition**  $DR_{p,\epsilon}$  is a discrepancy measure between  $\mu, \nu$   
 $\epsilon \in (0,1]$     $p \in (0,\infty)$     $d_{\mathcal{H}}$

$$DR_{p,\epsilon}(\mu, \nu) = \left( \int_{\epsilon}^1 d_{\mathcal{H}} \left( D_{\mu}^{\alpha}, D_{\nu}^{\alpha} \right)^p d\alpha \right)^{\frac{1}{p}}$$

## Properties

**Robustness**

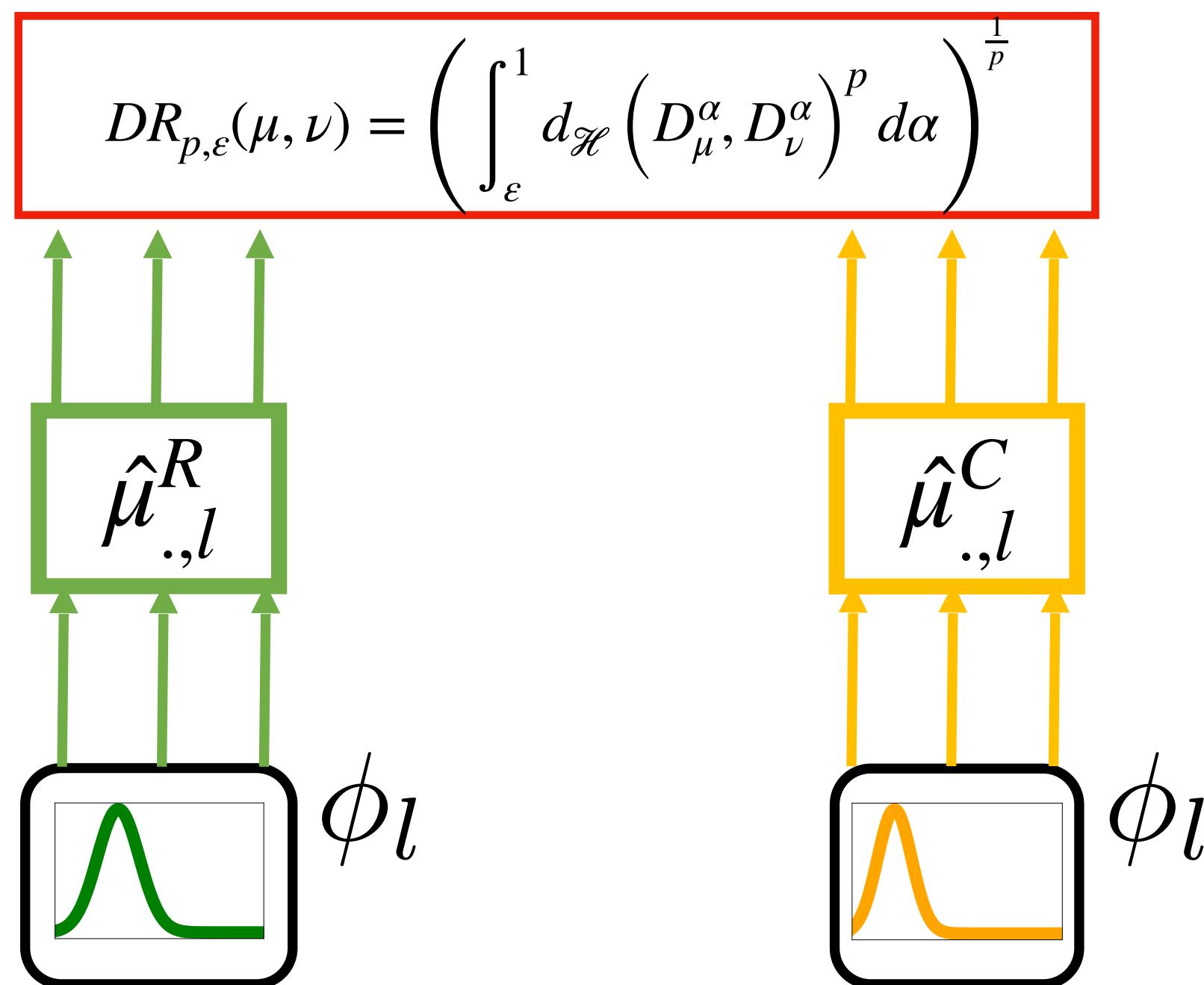


**Isometry invariance  
(translations,  
rotations)**

**Positive, symmetric,  
triangular inequality**

# DepthScore

G. Staerman, P. Mozharovskyi, P. Colombo, S. Clémencçon, F. d'Alché-Buc. A Pseudo-Metric between Probability Distributions based on Depth-Trimmed Regions.



R: The weather  
is cold today

C: It is freezing  
this morning

## Advantage

1. Deal with **paraphrases**
2. Include “**semantic**”

## Limitations

1. Use only **one layer**

# Results

---

## Experimental Setting

### Data2text Generation

- Results on **WebNLG 2020**
- **Correctness / Data Coverage / Relevance**
- Results on English only

### Summary Generation

- Results on **SummEval**
- Correlation with **pyramid score**
- Results on English only

## Summary Generation

	$DR_{p,\varepsilon}$	Abstractive			Extractive		
		$r$	$\tau$	$\rho$	$r$	$\tau$	$\rho$
Wasserstein	71.0	70.4	<u>71.1</u>	74.2	74.2	40.0	
Sliced-Wasserstein	70.1	68.7	71.0	72.4	73.9	<u>69.2</u>	
MMD	68.2	67.5	67.9	75.6	75.6	56.1	
BertScore	71.7	<u>71.9</u>	72.0	70.9	72.9	<b>73.8</b>	
MoverScore	<u>72.4</u>	<u>71.9</u>	<u>73.0</u>	<u>76.1</u>	<u>76.1</u>	47.4	
ROUGE-1	<b>73.5</b>	73.0	<b>74.4</b>	72.2	<u>74.0</u>	<u>69.1</u>	
ROUGE-2	73.0	<b>73.5</b>	73.0	55.1	53.2	<u>69.1</u>	
JS-2	68.9	6.8	69.8	<b>92.9</b>	5.5	19.0	

# Towards Multi-layers Metrics

---

**DepthScore & BertScore rely on a single layer !**

**Why are they design this way?**

**Before single layer embedding**

**Glove, Word2vec**

**Many prior work on single layer embeddings**

**WMD**

**What granularity is used ?**

**Word Level distance vs Sentence Level Distance**

**SentenceMover**

# BaryScore

Pierre Colombo, Guillaume Staerman, Chloé Clavel, Pablo Piantanida. Automatic Text Evaluation through the Lens of Wasserstein Barycenters. (oral) EMNLP 2021

# Optimal Transport

## Goal

Compute distance between probability measures  $(\mu, \nu)$

## Input Discret Measures

$$\nu = \sum_{j=1}^M \beta_j \delta_{x_j}$$

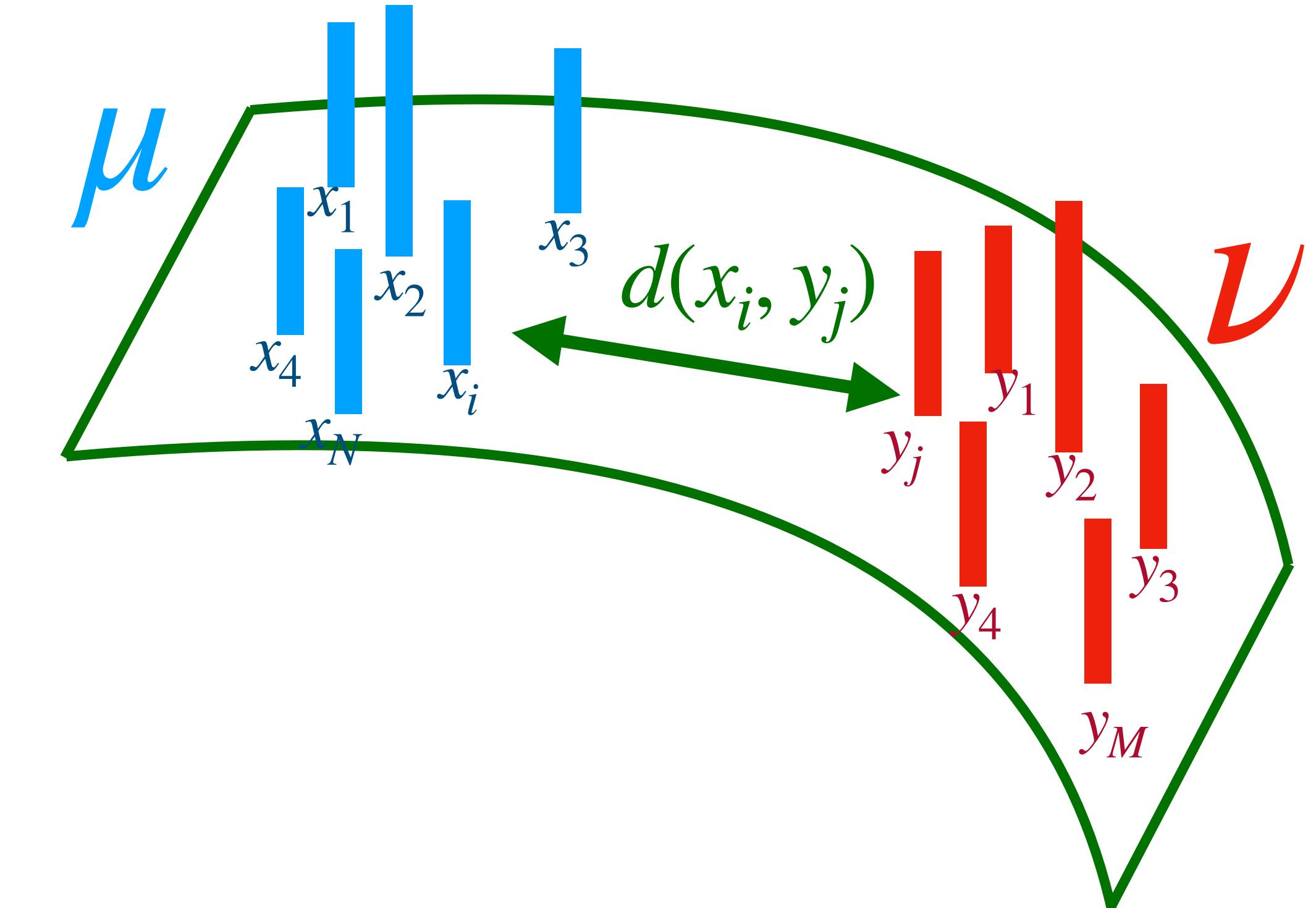
$$\mu = \sum_{i=1}^N \alpha_i \delta_{x_i}$$

## Cost Matrix

$$C = \begin{pmatrix} d(x_1, y_1) & \dots & d(x_1, y_M) \\ \dots & \dots & \dots \\ d(x_N, y_1) & \dots & d(x_N, y_M) \end{pmatrix}$$

## Transport Plan

$$\Pi = \begin{pmatrix} \pi_{11} & \dots & \pi_{1M} \\ \dots & \dots & \dots \\ \pi_{N1} & \dots & \pi_{NM} \end{pmatrix} \rightarrow \begin{matrix} \alpha_1 \\ \vdots \\ \alpha_N \end{matrix}, \begin{matrix} \beta_1 \\ \vdots \\ \beta_M \end{matrix}$$



## Wasserstein Distance

$$OT(\nu, \mu) = \min_{\Pi} \sum_{ij} C_{i,j} \times \Pi_{i,j}$$

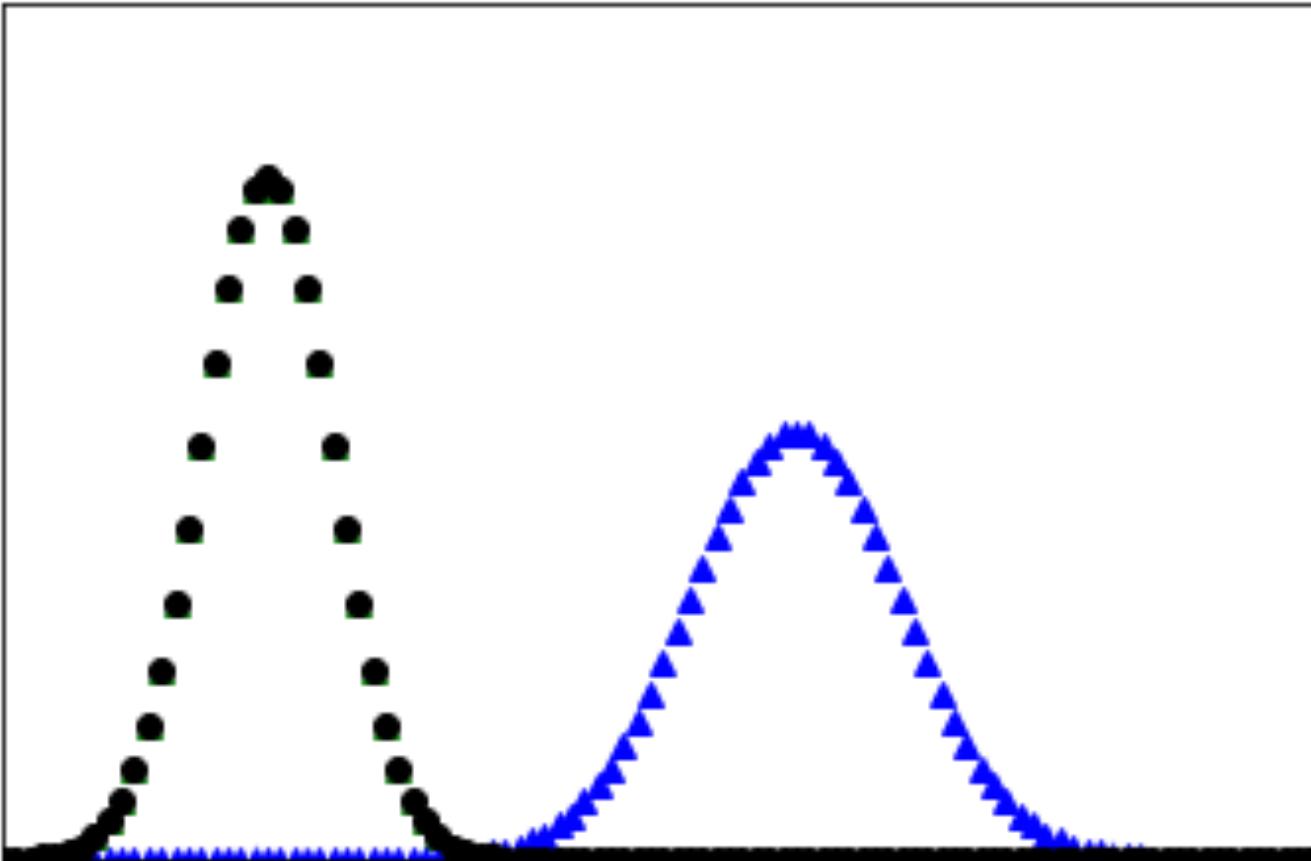
$$\Pi 1 = \alpha, \Pi^T 1 = \beta$$

# Wasserstein Barycenters

## Euclidean Interpolation

$$\nu = \sum_{i=1}^N \alpha_i l_2(\mu_i, \mu)$$

$$\nu = \alpha_i l_2(\mu_i, \mu) + (1 - \alpha_i) l_2(\mu_i, \mu)$$



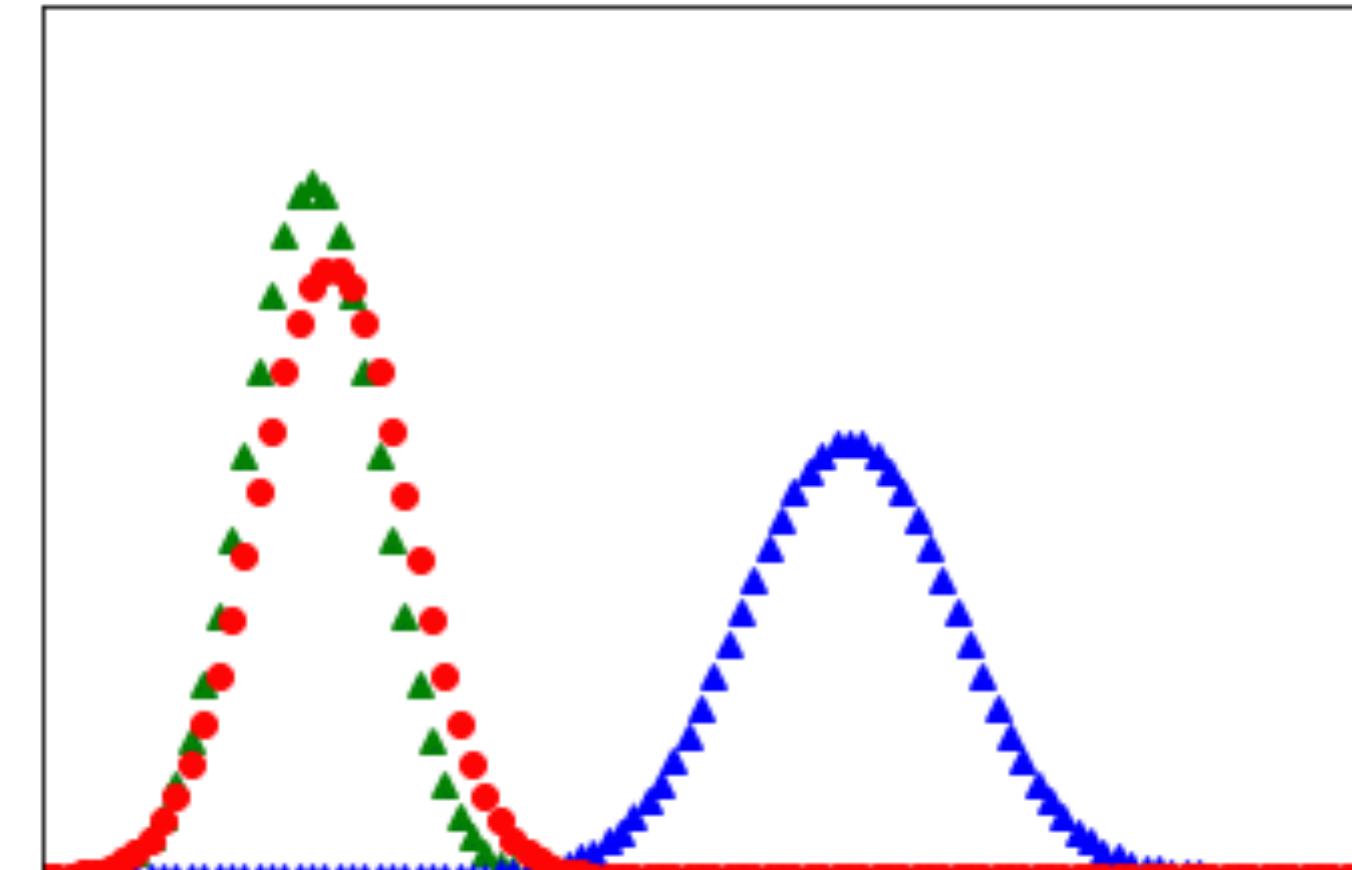
Do not look like a gaussian !

## Wasserstein Interpolation

$$\nu = \operatorname{argmin}_{\mu} \sum_{i=1}^N \alpha_i W(\mu_i, \mu)$$

$$\nu = \operatorname{argmin}_{\mu} \alpha_i W(\mu_i, \mu) + (1 - \alpha_i) W(\mu_i, \mu)$$

$\alpha_i$  Varies



Preserve the gaussian!

# BaryScore

**Reference:**  $R$

**Candidate:**  $C$

**Goal** : metric  $m : (\boxed{C}, \boxed{R}) \mapsto m(C, R) \in \mathbb{R}_+$

## Algorithm

1. Find the Wasserstein barycentric distributions of BERT layers for  $C$  and  $R$

2. Evaluate these barycentric distributions using the Wasserstein distance.

---

### Algorithm 1 BaryScore

---

**INPUT:**  $C = \{\omega_1^c, \dots, \omega_{n_c}^c\}$ ,  $R = \{\omega_1^r, \dots, \omega_{n_r}^r\}$ ,  
 $(\phi_1, \dots, \phi_L)$  pre-trained layers from BERT or ELMo.

**Compute layers embeddings:**

$\phi_\ell(C)$  and  $\phi_\ell(R)$  for every  $1 \leq \ell \leq L$ .

**Compute measures:**  $\{\hat{\mu}_{C,\ell}, \hat{\mu}_{R,\ell}\}_{\ell=1}^L$ .

**Compute Wasserstein barycenters:**

$$\hat{\mu}_C = \operatorname{argmin}_{\hat{\mu}} \sum_{\ell=1}^L \mathcal{W}(\hat{\mu}_{C,\ell}, \hat{\mu}),$$

$$\hat{\mu}_R = \operatorname{argmin}_{\hat{\mu}} \sum_{\ell=1}^L \mathcal{W}(\hat{\mu}_{R,\ell}, \hat{\mu}),$$

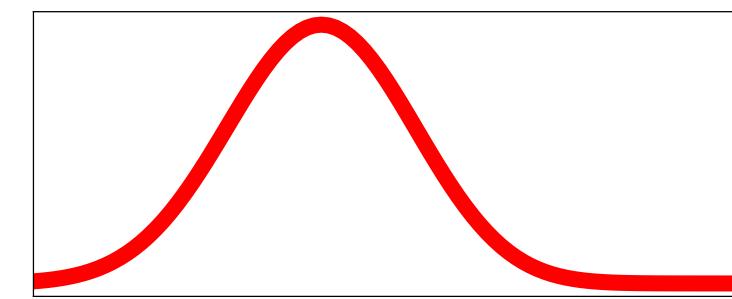
---

**OUTPUT:**  $\mathcal{W}(\hat{\mu}_R, \hat{\mu}_C)$ .

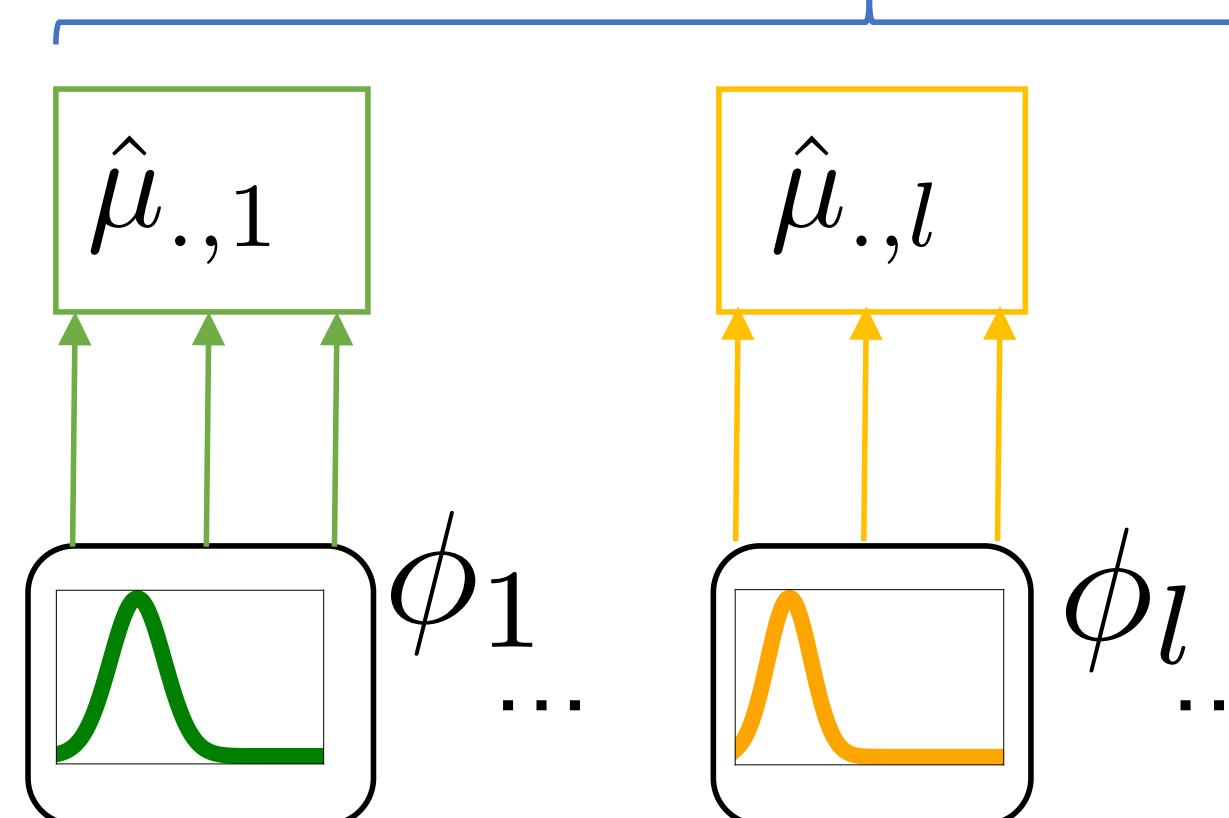
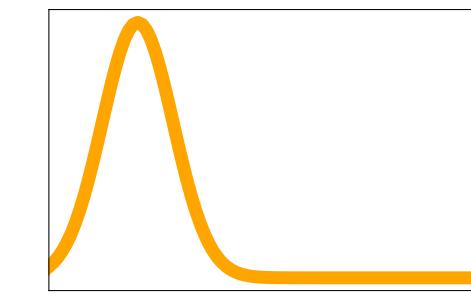
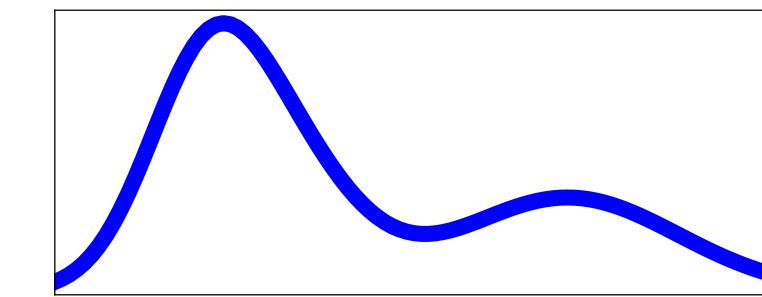
---

# BaryScore vs BertScore vs MoverScore

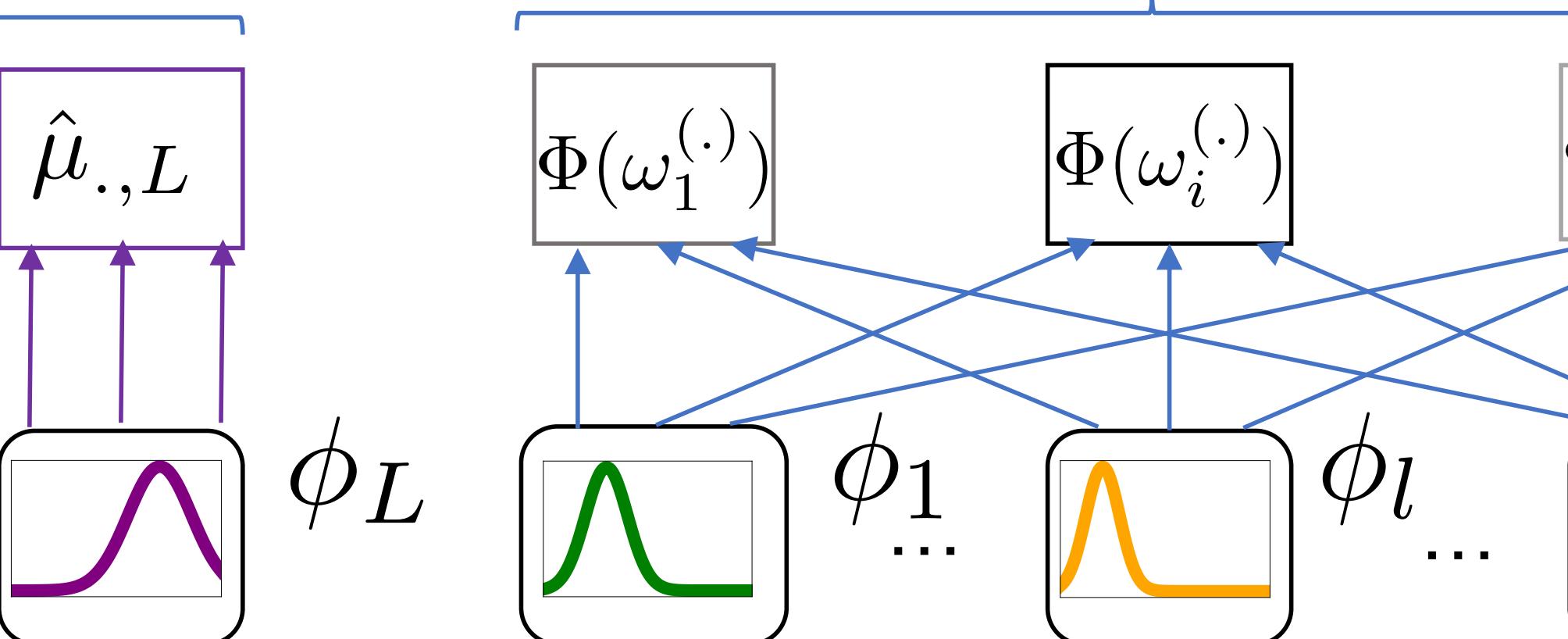
## Comparison between aggregation functions



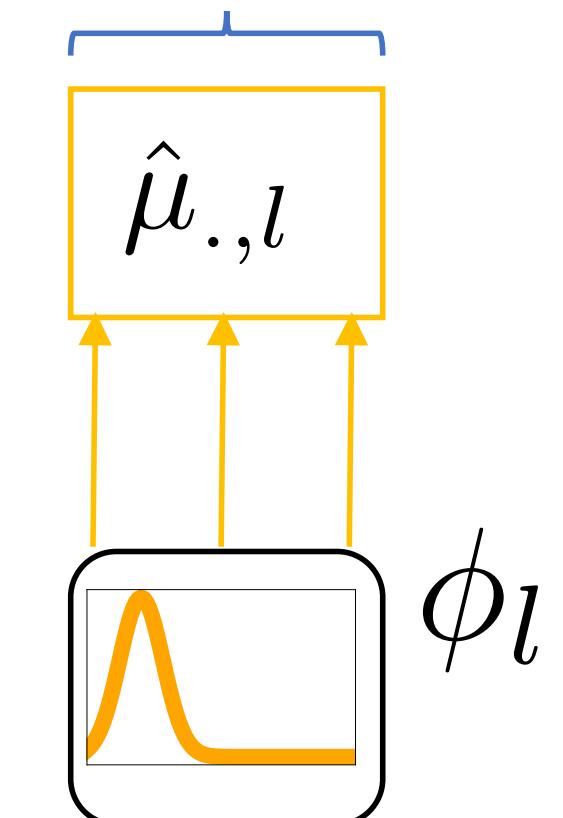
$$\hat{\mu}_{\cdot} = \operatorname{argmin}_{\hat{\mu}} \sum_{\ell=1}^L \mathcal{W}(\hat{\mu}_{\cdot,\ell}, \hat{\mu})$$



(a)



(b)



(c)

**BaryScore**

**MoverScore**

**BertScore**

## Experimental Setting

---

### Machine Translation

- Results on **WMT17/WMT18**
- All metrics are measures on en only
- Pairs includes cs-en de-en ru-en fi-en ro-en tr-en

### Data2text Generation

- Results on **WebNLG 2020**
- Correctness / Data Coverage / Relevance
- Results on English only

### Image Captioning

- Results on **MSCOCO**
- Results on English only

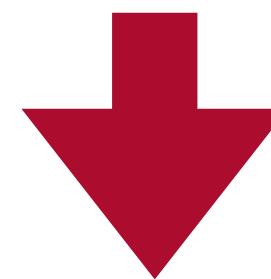
### Summary Generation

- Results on **SummEval**
- Correlation with **pyramid score**
- Results on English only

# Results

## Task

(John\ Blaha birthDate 1942\ 08\ 26)  
(John\ Blaha birthPlace San\ Antonio)  
(John\ E\ Blaha job Pilot)



John Blaha, born in San  
Antonio on 1942-08-26,  
worked as a pilot

## Criterion

Correctness

Data coverage

Relevance

Metric	Correctness			Data Coverage			Relevance		
	$r$	$\rho$	$\tau$	$r$	$\rho$	$\tau$	$r$	$\rho$	$\tau$
Correct	100.0	100.0	100.0	97.6	85.2	73.3	99.1	89.7	75.0
DataC	85.2	97.6	73.3	100.0	100.0	100.0	96.0	93.8	81.6
Relev	89.7	99.1	75.0	96.0	93.8	81.6	100.0	100.0	100.0
BaryS	<b>91.7</b>	<b>90.0</b>	<b>78.3</b>	<b>87.8</b>	78.2	61.6	<b>89.4</b>	<b>82.6</b>	70.0
BaryS <sup>+</sup>	90.5	89.5	76.6	87.7	<b>85.0</b>	<b>70.0</b>	89.2	86.4	71.6
BertS	85.5	85.4	75.5	74.1	68.2	55.5	85.5	79.4	65.0
MoverS	84.1	84.1	73.3	78.7	66.2	53.3	82.1	77.4	65.0
BLEU	77.6	66.3	60.0	55.7	50.2	36.6	63.0	65.2	51.6
R-1	80.6	65.0	65.0	76.5	76.3	60.3	64.3	69.2	56.7
R-2	73.6	63.3	58.3	54.7	43.1	35.0	62.0	60.8	46.7
R-WE	60.9	73.4	60.0	40.2	58.2	40.1	49.9	64.1	48.3
METEOR	86.5	66.3	70.0	77.3	50.2	46.6	82.1	65.2	58.6
TER	79.6	78.3	58.0	69.7	58.2	38.0	75.0	70.2	<b>77.6</b>

Correlation score for different coefficient Pearson  $r$ , Spearman  $\rho$  and Kendall  $\tau$

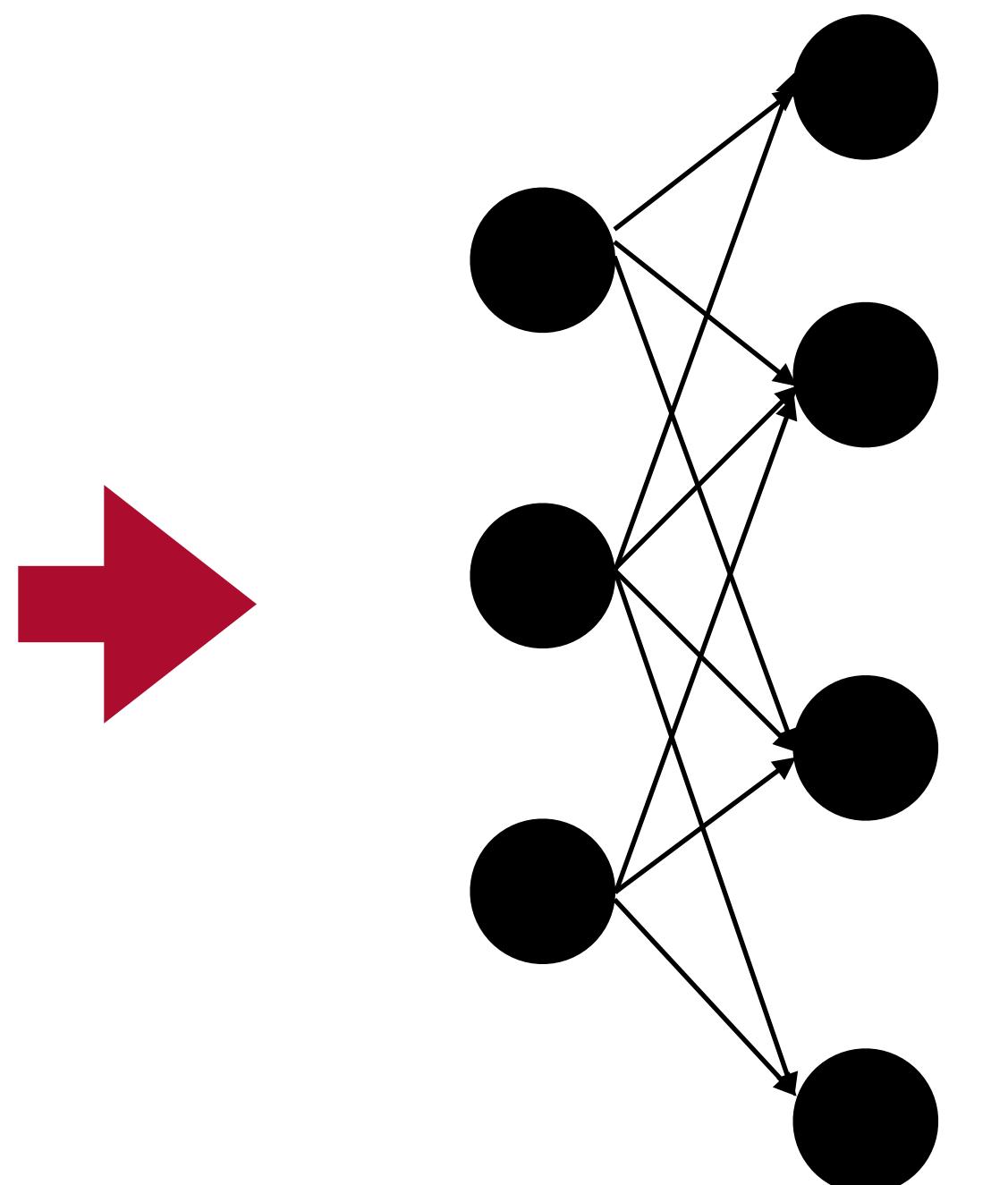
# InfoLM

Pierre Colombo, Chloé Clavel and Pablo Piantanida. InfoLM: A New Metric to Evaluate Summarization & Data2Text Generation. AAAI 2022

# Statistical Measures of Similarity

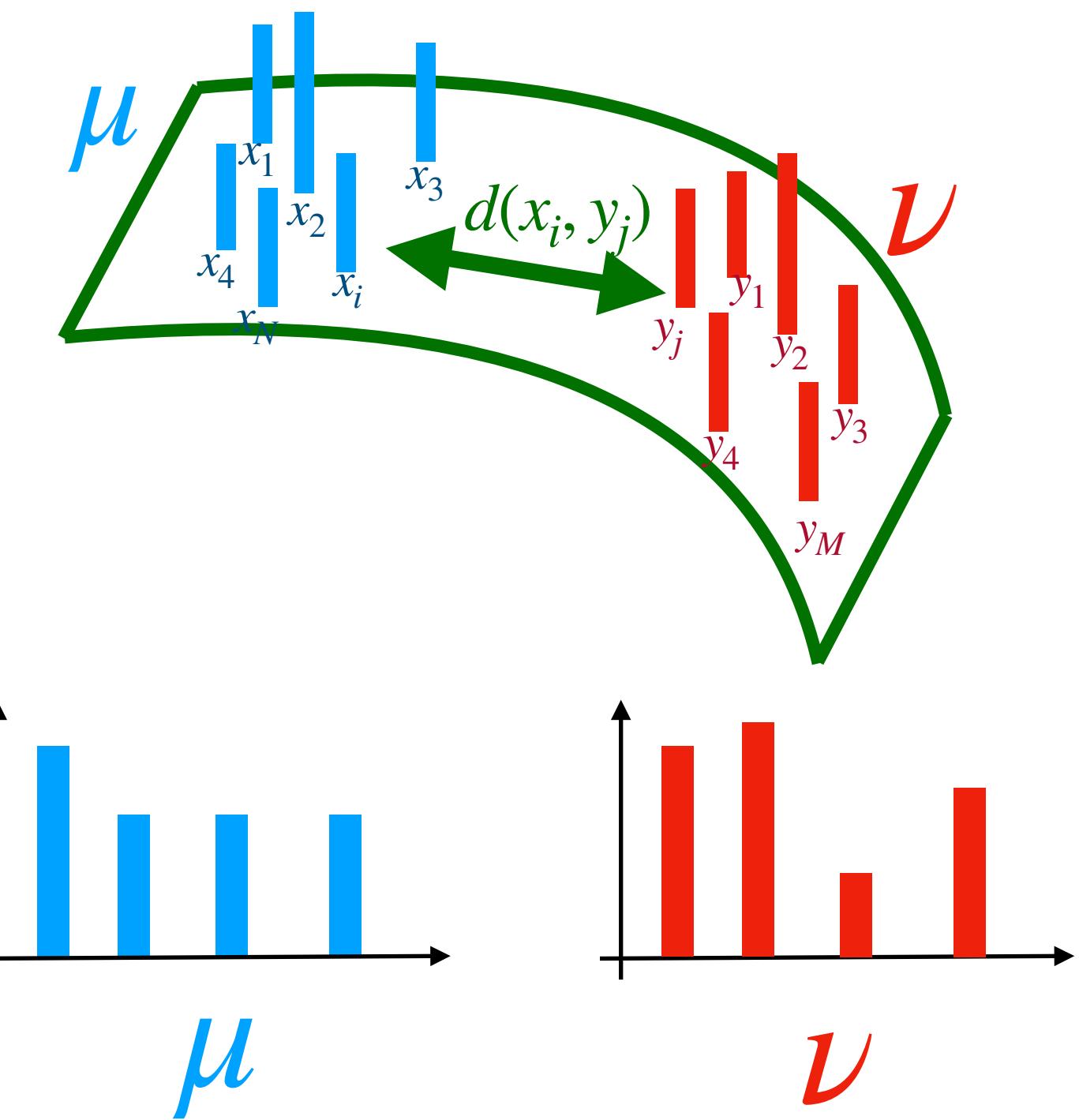
Hello, Chicago.  
If there is anyone out  
there who still doubts  
that America is a place  
where all things are  
possible, who still  
wonders if the dream of  
our founders is alive in  
our time, [...].  
Yes we can!

Input Text



Neural Network

High dimensional data



Soft Probabilities

# Existing Methods

# InfoLM

## Edit Based

Snover et al. 2006

### Operations

- Insertion (I)
- Deletion (D)
- Substitution (S).

tailor -> sailor (**S**)

sailor -> sailir (**S**)

sailir -> sailin (**S**)

sailin -> sailing (**I**)

Distance is 4 !

## N-gram Based

Papineni et al. 2002

C : I like these very nice pies !

R : I like those cakes !

### Unigrams

C : I like these very nice pies !

R : I like those cakes !

### Bigrams

C : I like these very nice pies !

R : I like those cakes !

## Embedding Based

### Word Mover distance

Kusner et al. 2015

### BertScore

Zhang et al. 2019

### MoverScore

Zhao et al. 2019

### Sentence Mover

Clark et al. 2019

# Assumptions for InfoLM

---

**Goal** Compute a similarity score between R and C.

**Tools** Use a pretrained MLM

MLM predicts a distribution over  $\Omega$

$$p_{\Omega}(\cdot | [R]^i)$$

Use a measure of information

$$\mathcal{I} : [0,1]^{|\Omega|} \times [0,1]^{|\Omega|}$$

Name	Notation	Domain	Expression
$\alpha$ -divergence (Csiszár 1967)	$\mathcal{D}_{\alpha}$	$\alpha \notin \{0, 1\}$	$\frac{1}{\alpha(\alpha-1)}(1 - \sum q_i^{1-\alpha} p_i^{\alpha})$
$\gamma$ divergence (Fujisawa and Eguchi 2008)	$\mathcal{D}_{\gamma}^{\beta}$	$\beta \notin \{0, -1\}$	$\frac{1}{\beta(\beta+1)} \log \sum p_i^{\beta+1} + \frac{1}{\beta+1} \log \sum q_i^{\beta+1} - \frac{1}{\beta} \log \sum p_i q_i^{\beta}$
AB Divergence (Cichocki, Cruces, and Amari 2011)	$\mathcal{D}_{sAB}^{\alpha,\beta}$	$(\alpha, \beta) \in (\mathbb{R}^*)^2$ $\beta + \alpha \neq 0$	$\frac{1}{\beta(\beta+\alpha)} \log \sum p_i^{\beta+\alpha} + \frac{1}{\beta+\alpha} \log \sum q_i^{\beta+\alpha} - \frac{1}{\beta} \log \sum p_i^{\alpha} q_i^{\beta}$
$\mathcal{L}_1$ distance	$\mathcal{L}_1$		$\sum  p_i - q_i $
$\mathcal{L}_2$ distance	$\mathcal{L}_2$		$\sqrt{\sum (p_i - q_i)^2}$
$\mathcal{L}_{\infty}$ distance	$\mathcal{L}_{\infty}$		$\max_i  p_i - q_i $
Fisher-Rao distance	$R$		$\frac{2}{\pi} \arccos \sum \sqrt{p_i \times q_i}$

# Intuition of InfoLM

**Goal** Compute a similarity score between R and C.

**Equivalence for masked contexts**

$$\mathcal{I} : [0,1]^{|\Omega|} \times [0,1]^{|\Omega|}$$

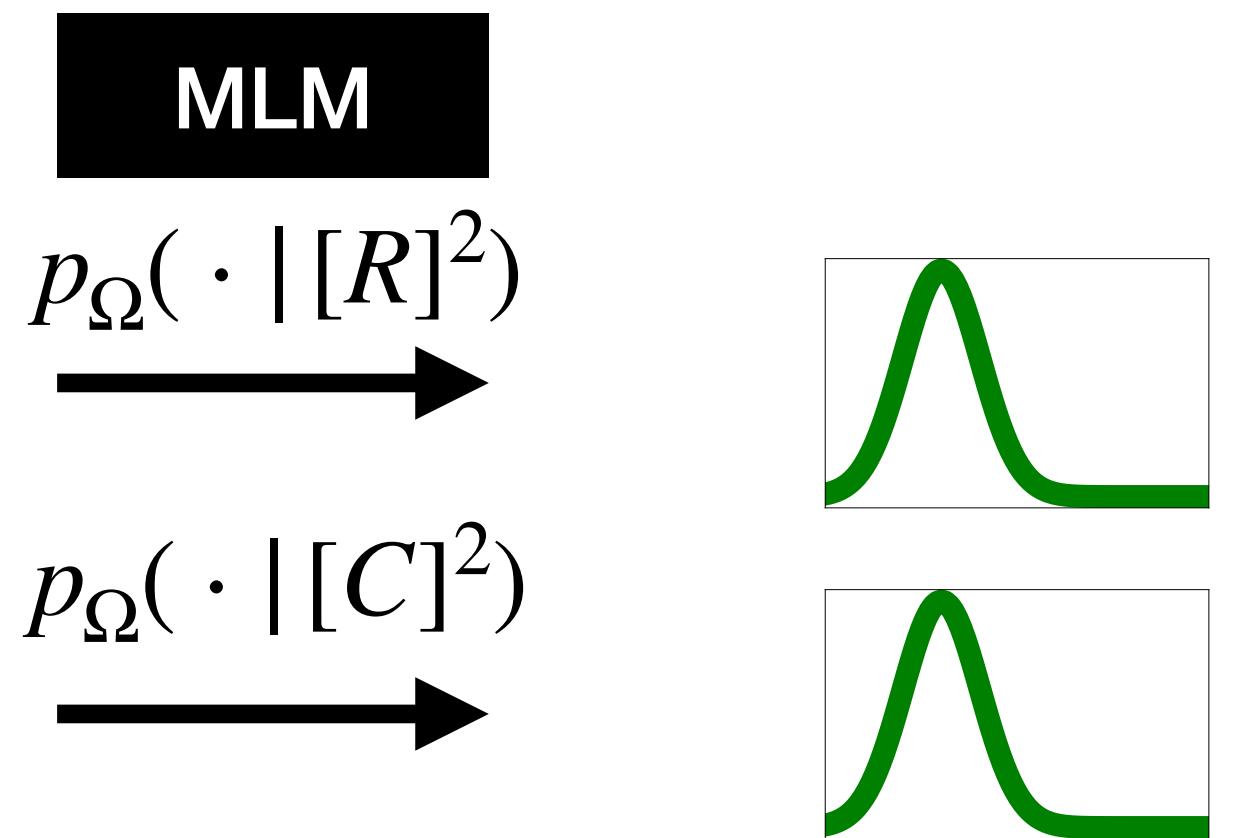
MLM predicts a distribution over  $\Omega$

$$p_{\Omega}(\cdot | [R]^i)$$

**Similar context**

R: It is [MASK] today.

C: It is [MASK] this morning !

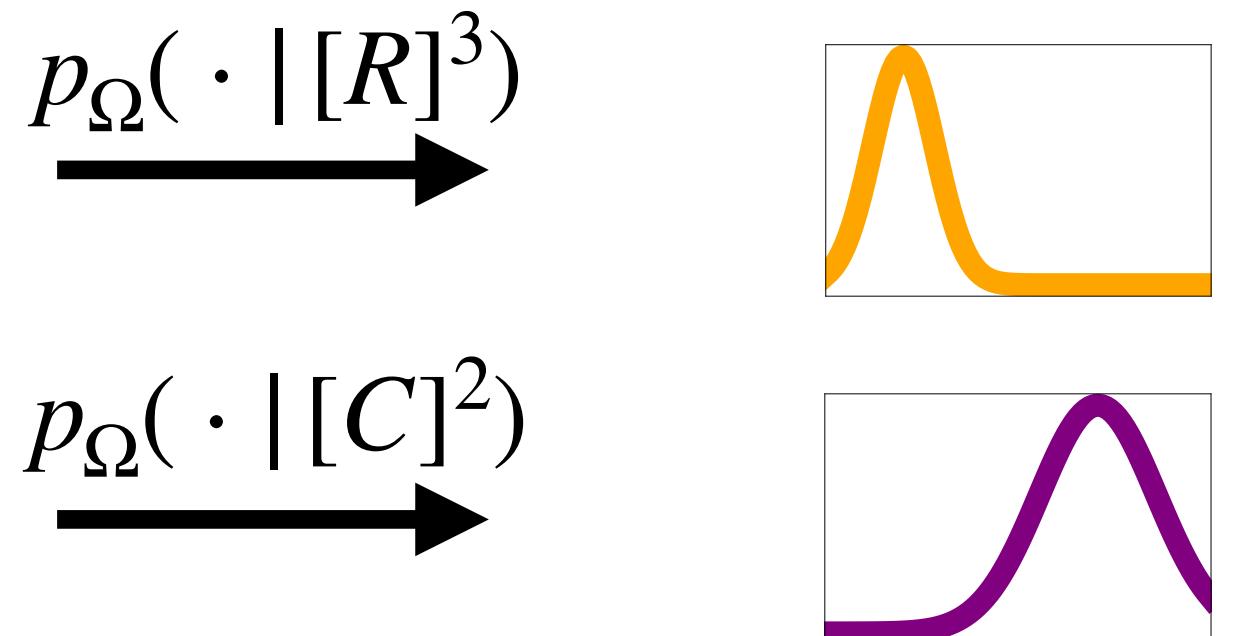


$$\mathcal{I}(p_{\Omega}(\cdot | [R]^2), p_{\Omega}(\cdot | [C]^2)) \sim 0$$

**Dissimilar context**

R: It is cold [MASK]

C: It is [MASK] this morning !

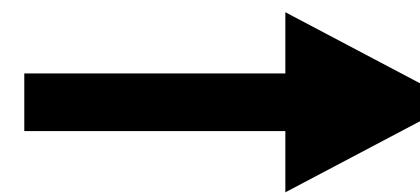


$$\mathcal{I}(p_{\Omega}(\cdot | [R]^3), p_{\Omega}(\cdot | [C]^2)) \gg 0$$

# Context Aggregation

**Goal** Compute a similarity score between R and C.

How to aggregate contexts?



Weighted Sum!

Reference

[MASK] is cold today.

It is [MASK] today.

...

It is cold today [MASK]

$$P \triangleq \frac{1}{5} \sum_{k=0}^4 \gamma_k \times p_{\Omega}(\cdot | [R]^k)$$

$$\text{InfoLM}(R, C) \triangleq \mathcal{J}(P, Q)$$

Candidate

[MASK] is freezing this morning !

It is [MASK] this morning !

...

It is freezing this morning [MASK]

$$Q \triangleq \frac{1}{6} \sum_{k=0}^5 \gamma_k \times p_{\Omega}(\cdot | [C]^k)$$

## Experimental Setting

---

### Data2text Generation

- Results on **WebNLG 2020**

Gardent et al. 2017

- **Correctness / Data Coverage / Relevance**  
**Fluency / Text Structure**

Ferreira et al. (2020)

Perez-Beltrachini et al 2016

- Results on English only

### Summary Generation

- Results on **SummEval**

Nallapati et al. 2016)  
Bhandari et al. (2020)

- **Correlation with pyramid score**

Nenkova and Passonneau 2004

- Results on English only

# Results

---

Task

(John\\_Blaha birthDate 1942\\_08\\_26)  
 (John\\_Blaha birthPlace San\\_Antonio)  
 (John\\_E\\_Blaha job Pilot)



John Blaha, born in San Antonio on 1942-08-26, worked as a pilot

Metric	Correctness			Data Coverage			Fluency			Relevance			Text Structure		
	r	$\rho$	$\tau$	r	$\rho$	$\tau$	r	$\rho$	$\tau$	r	$\rho$	$\tau$	r	$\rho$	$\tau$
Correct	100.0	100.0	100.0	97.6	85.2	73.3	80.0	81.1	61.6	99.1	89.7	75.0	80.1	80.8	60.0
DataC	85.2	97.6	73.3	100.0	100.0	100.0	71.8	51.7	38.3	96.0	93.8	81.6	71.6	51.4	36.6
Fluency	81.1	80.0	61.6	71.8	51.7	38.3	100.0	100.0	100.0	77.0	61.4	46.6	99.5	99.7	98.3
Relev	89.7	99.1	75.0	96.0	93.8	81.6	77.0	61.4	46.6	100.0	100.0	100.0	77.2	61.1	45.0
TextS	80.8	80.1	60.0	71.6	51.4	36.6	99.5	99.7	98.3	77.2	61.1	45.0	100.0	100.0	100.0
$\mathcal{D}_{AB}$	88.8	<u>89.3</u>	<u>76.6</u>	<u>81.8</u>	<u>82.6</u>	<u>70.0</u>	86.6	92.0	76.6	<u>89.8</u>	<u>87.9</u>	<u>73.3</u>	86.6	91.4	75.0
$\mathcal{D}_\alpha$	88.8	<u>89.3</u>	<u>76.6</u>	<u>81.8</u>	<u>82.6</u>	<u>70.0</u>	86.6	92.0	76.6	<u>89.8</u>	<u>87.9</u>	<u>73.3</u>	86.6	91.4	75.0
$\mathcal{D}_\beta$	81.4	50.0	71.6	48.4	79.7	65.0	44.8	84.7	76.6	49.3	72.3	60.0	48.0	83.8	75.0
$\mathcal{L}_1$	75.2	33.8	61.6	32.4	53.8	40.0	22.7	83.5	73.3	32.2	57.9	45.0	25.6	83.2	71.6
$\mathcal{R}$	<u>89.7</u>	86.0	75.0	78.7	70.5	51.6	<u>93.3</u>	<u>95.7</u>	<u>85.3</u>	87.6	84.4	70.0	<u>92.4</u>	<u>93.8</u>	<u>81.6</u>
JS	79.4	81.1	70.0	69.3	75.5	60.0	89.4	91.4	75.0	81.7	70.5	60.0	91.9	91.1	73.3
BertS	<u>85.5</u>	83.4	<u>73.3</u>	74.7	<u>68.2</u>	53.3	<u>92.3</u>	<u>95.5</u>	<u>85.0</u>	<u>83.3</u>	<u>79.4</u>	<u>65.0</u>	<u>91.9</u>	<u>95.0</u>	<u>83.3</u>
MoverS	84.1	<u>84.1</u>	<u>73.3</u>	<u>78.7</u>	66.2	<u>53.3</u>	91.2	92.1	78.3	82.1	77.4	65.0	90.1	91.4	76.3
BLEU	77.6	66.3	60.0	55.7	50.2	36.6	<u>89.4</u>	90.5	78.3	63.0	65.2	51.6	88.5	89.1	76.6
R-1	80.6	65.0	65.0	61.1	<u>59.6</u>	48.3	76.5	76.3	60.3	64.3	<u>69.2</u>	56.7	75.9	77.5	58.3
METEOR	<u>86.5</u>	<u>66.3</u>	<u>70.0</u>	<u>77.3</u>	50.2	46.6	86.7	90.5	78.3	<u>82.1</u>	<u>65.2</u>	58.6	86.2	89.1	76.6
TER	79.6	78.3	58.0	69.7	58.2	38.0	89.1	<u>93.5</u>	<u>80.0</u>	75.0	70.2	<u>77.6</u>	89.5	91.1	78.6

# Summary

---

**Summary** In this presentation we have explored the problem of NLG evaluation with a specific focus on reference based evaluation

DepthScore

New Metrics

BaryScore

InfoLM

**Future Work** Most of SOTA metrics are based on pretrained representations and lack of interpretability.

Introduce confidence scores.

Build news datasets.

# Further Readings

---

**Many Future works can be drawn from the works I presented !**

## Papers

**Automatic Machine Translation  
Evaluation in Many Languages via  
Zero-Shot Paraphrasing**

**Uncertainty-Aware Machine  
Translation Evaluation**

**Pushing the Right Buttons:  
Adversarial Evaluation of  
Quality Estimation**

**Better than Average: Paired  
Evaluation of NLP Systems**

## Researchers

**Matt Post**  
<https://scholar.google.com/citations?user=4w7LhxsAAAAJ&hl=en>

**Andre Martins**  
<https://scholar.google.com/citations?user=mT7ppvwAAAAJ&hl=en>

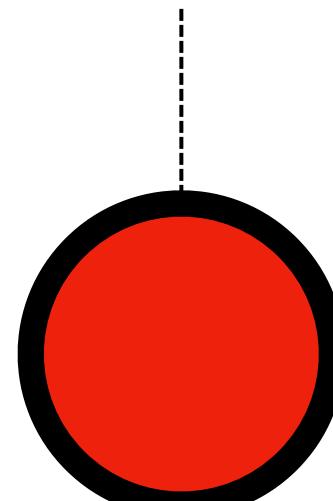
**Lucia Specia**  
[https://scholar.google.co.uk/citations?user=wVI\\_z8kAAAAJ&hl=en](https://scholar.google.co.uk/citations?user=wVI_z8kAAAAJ&hl=en)

**Maxime Peyrard**  
<https://scholar.google.fr/citations?user=RFMdKLMAAAJ&hl=en>

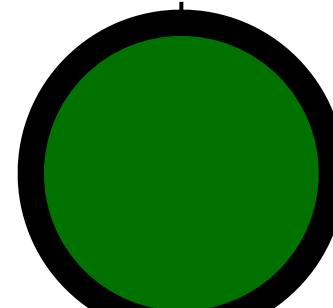
# **Other lines of research I am pursuing and you could be interested in!**

---

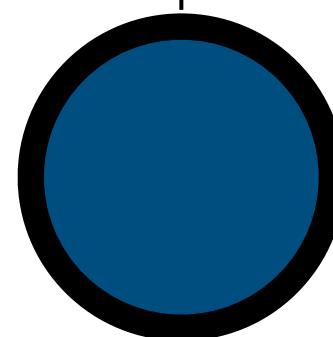
## **Machine Learning & Security**



**Learning Disentangled Representations**

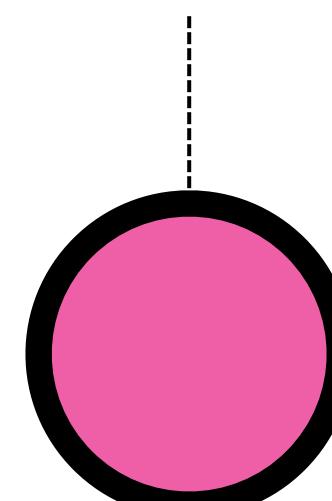


**Out of Distribution detection**

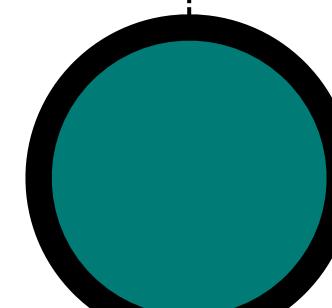


**Miss-classification detection**

## **Dialog Systems**



**Learning Representations of Dialogs**



**Multimodal Learning**

**Check out my scholar:**



## Acknowledgements

---

Computing Grant from GENCI.



**This work would no have been possible without my co-authors**

**Guillaume Staerman, Chloé Clavel, Pavlo Mozharovskyi, Stephan Cléménçon, Florence d'Alché-Buc,  
Pablo Piantanida**



If you want to contact me:

[colombo.pierre@gmail.com](mailto:colombo.pierre@gmail.com)

Website

Twitter

Scholar

# **Automatique Evaluation of Natural Language Generation using Measures of Similarity**

**Thanks for your attention!**

# Conversational Agents

## Definition

A conversational agent is any dialogue system that not only conducts natural language processing but also responds automatically using human language.



## Natural Language Understanding (NLU)

The conversation agent has to understand the user.

## Natural Language Generation (NLG)

The conversation agent has to produce a response.



## Existing Methods

### Edit Based

Snover et al. 2006

#### Operations

- Insertion (I)
- Deletion (D)
- Substitution (S).

tailor -> sailor (**S**)

sailor -> sailir (**S**)

sailir -> sailin (**S**)

sailin\_ -> sailing (**I**)

Distance is 4 !

### N-gram Based

Papineni et al. 2002

C : I like these very nice pies !

R : I like those cakes !

#### Unigrams

C : I like these very nice pies !

R : I like those cakes !

#### Bigrams

C : I like these very nice pies !

R : I like those cakes !

## MoverScore & BaryScore

### Embedding Based

#### Word Mover distance

Kusner et al. 2015

#### BertScore

Zhang et al. 2019

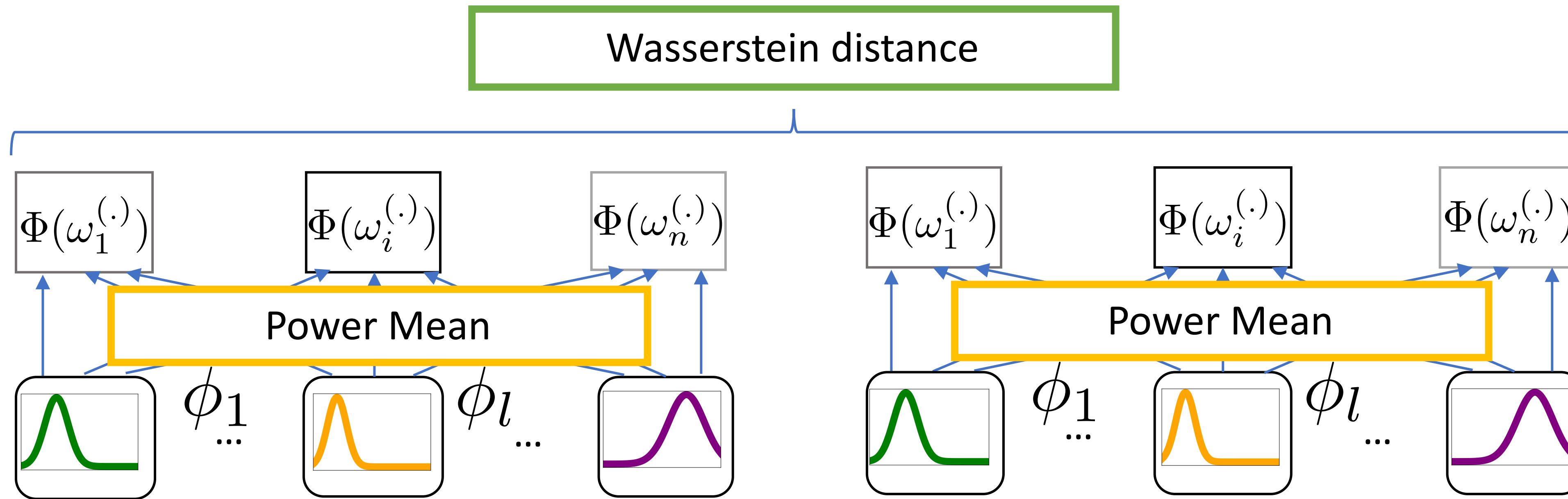
#### MoverScore

Zhao et al. 2019

#### Sentence Mover

Clark et al. 2019

# MoverScore



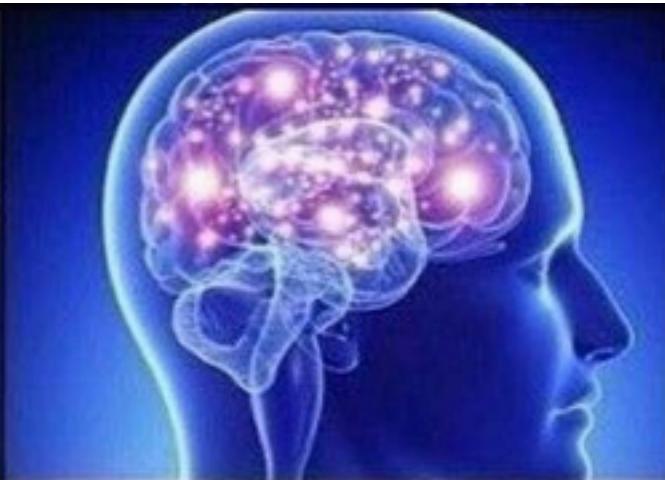
## Advantage

1. Deal with paraphrases
2. Include “semantic”
3. Use several layers

## Limitations

1. Use arbitrary sequence of operation  
(euclidean aggregation function  
Wasserstein distance)

# Towards BaryScore



## A novel metric called BaryScore

Previously

1 Take one layer

2 Do a series of operations  
(Wasserstein)

BertScore

DepthScore

1 Take several layers

2 Aggregate using  
euclidean distance

3 Do a series of  
Operations (Wasserstein)

MoverScore

Best of all worlds

1 Take several layers

2 Aggregate using  
Wasserstein distance

3 Do a series of  
operation (Wasserstein)

BaryScore