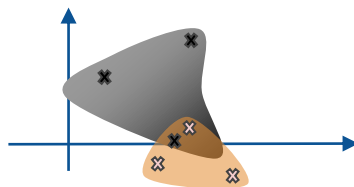


# Analysis and prediction of NLP models via task embeddings



Damien Sileo and Marie-Francine Moens  
KU Leuven

[damien.sileo@kuleuven.be](mailto:damien.sileo@kuleuven.be)



Horizon 2020  
Advanced Grant 788506



# Introduction

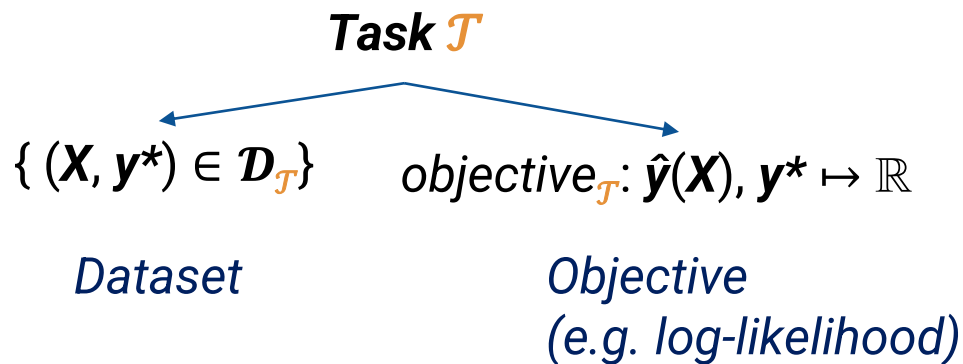
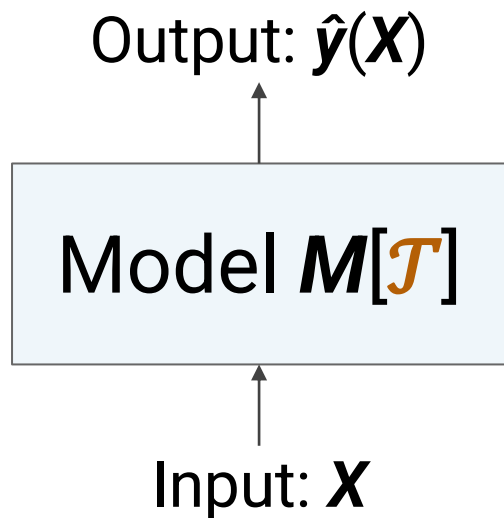
---

Transfer & task embeddings

# Contributions

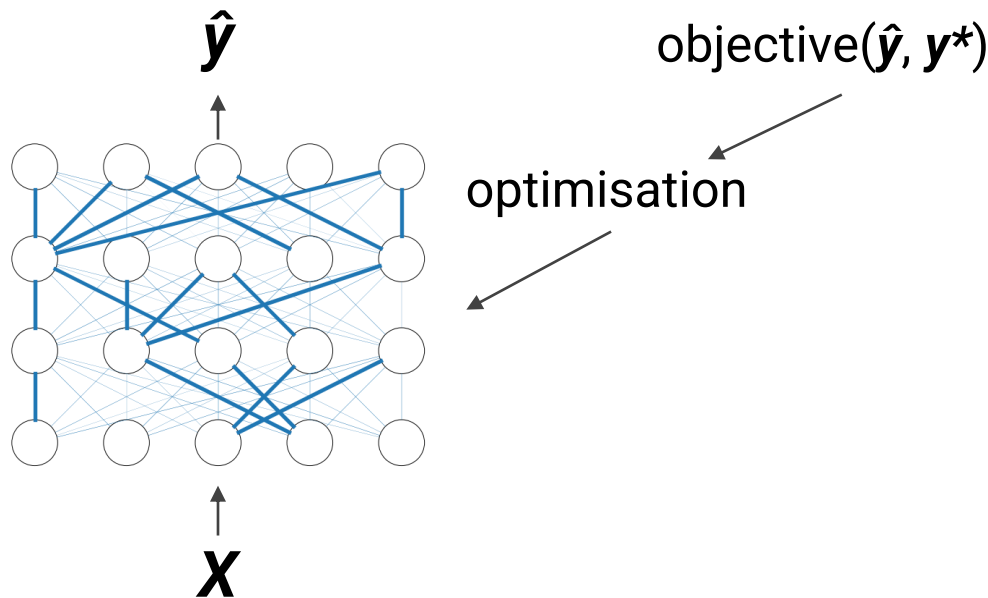
New resource: metaEval  
Experimental setting  
Analyses of task embeddings

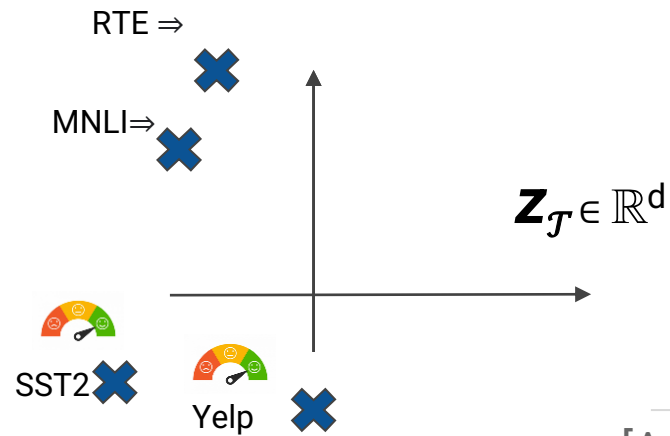




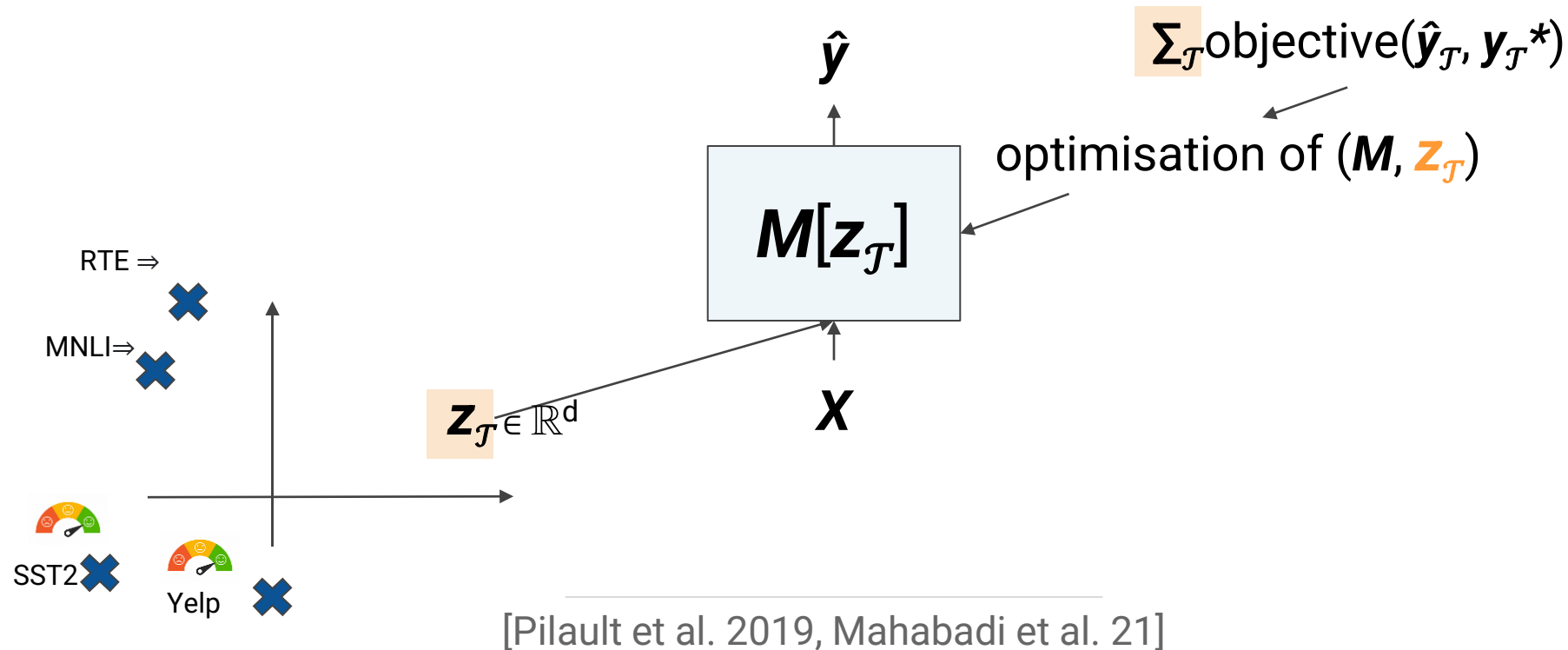
$M[\emptyset \rightarrow \text{RTE}]$ 
 $\prec$ 
 $M[\emptyset \rightarrow \text{MNLI} \rightarrow \text{RTE}]$   
 SOURCE TARGET

transfer





[Achille et al. 2019, Vu et al. 2019]



## Introduction

---

Transfer & task embeddings

## Contributions

New resource: **metaEval**

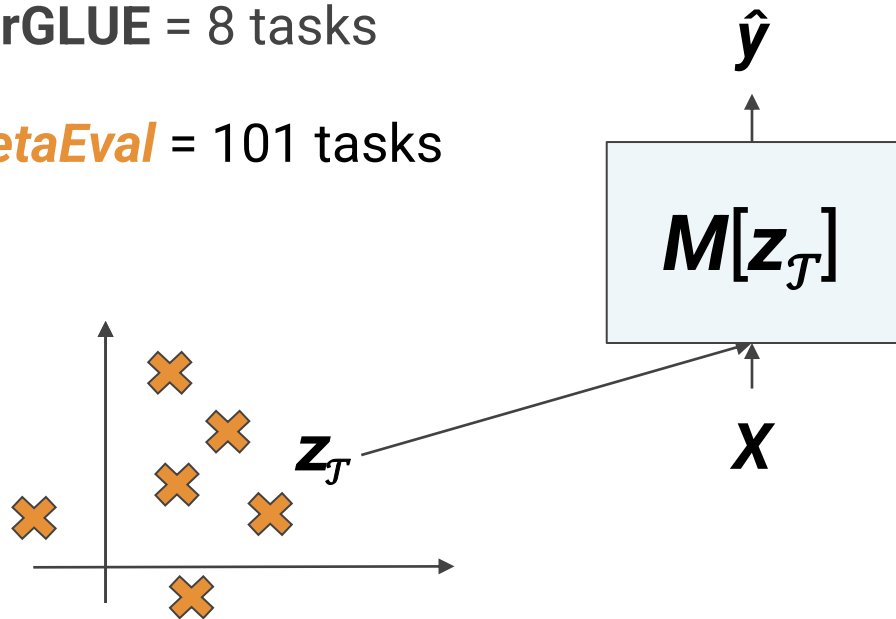
Experimental setting

Analyses of task embeddings

GLUE = 9 tasks

SuperGLUE = 8 tasks

*MetaEval* = 101 tasks





Huggingface Datasets (Wolf et al. 20)

GLUE (Wang et al. 19)

SuperGLUE (Wang et al. 19)

Other classification tasks

+ pragmEval (Sileo et al. 21)

+ TweetEval (Barbieri et al. 20)

+ Recast (Poliak et al. 18)

+ Blimp (Warstadt et al. 20)

+ Ethics (Hendrycks et al. 21)

+ CrowdFlowers (Van pelt et al. 12)

## 101 tasks

Unified 🌐 dataset format: (train/val/test)

sentence1	sentence2	label
sentence	-	label

```
dataset = load_and_align('health_fact')
```

<https://github.com/sileod/metaeval>

## Introduction

---

Transfer & task embeddings

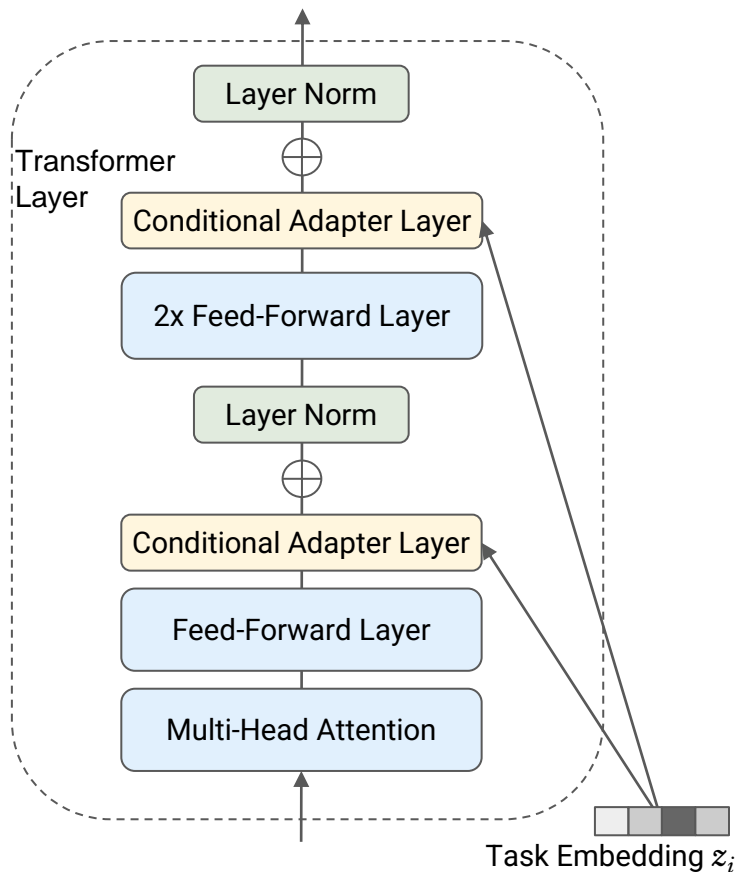
## Contributions

New resource: metaEval

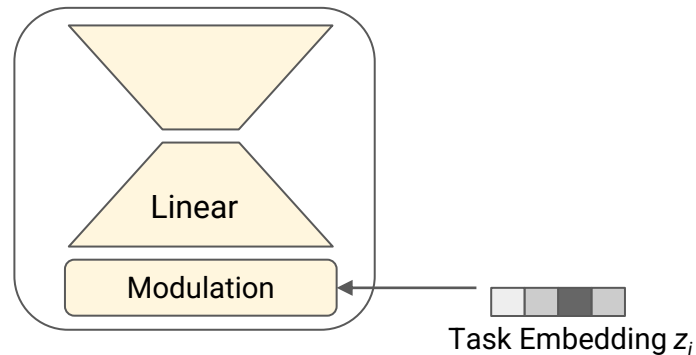
**Experimental setting**

Analyses of task embeddings

# Conditional Adapters



## Conditional Adapter Layer



[Pilault et al. 21]

trained

frozen

Fine-tuning method	MetaEval test accuracy	Trained encoder parameters	Task-specific trained encoder parameters
Majority Class	42.9	-	-
Full Fine-Tuning (1 model/task)	76.9	124M	124M
Adapter (1 adapter/task)	67.8	10M	10M
Conditional Adapter (Mahabadi et al. 2021)	75.6	38M	512
Conditional Adapter (Pilault et al. 2021)	<b>79.7</b>	10M	32
z=TextEmb task embedding (Vu et al. 2020)	69.9	10M	
z=Fisher-information task embedding (Fisher et al 2020)	67.5	10M	

## Introduction

---

Transfer & task embeddings

## Contributions

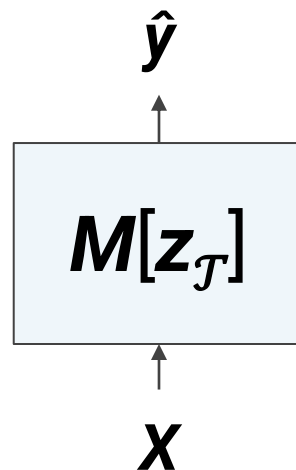
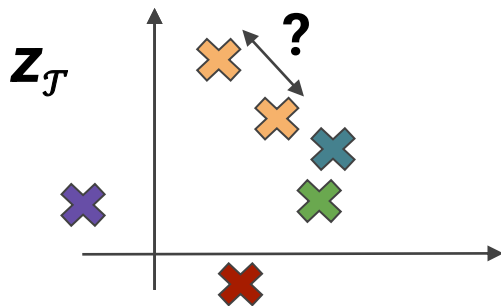
New resource: metaEval

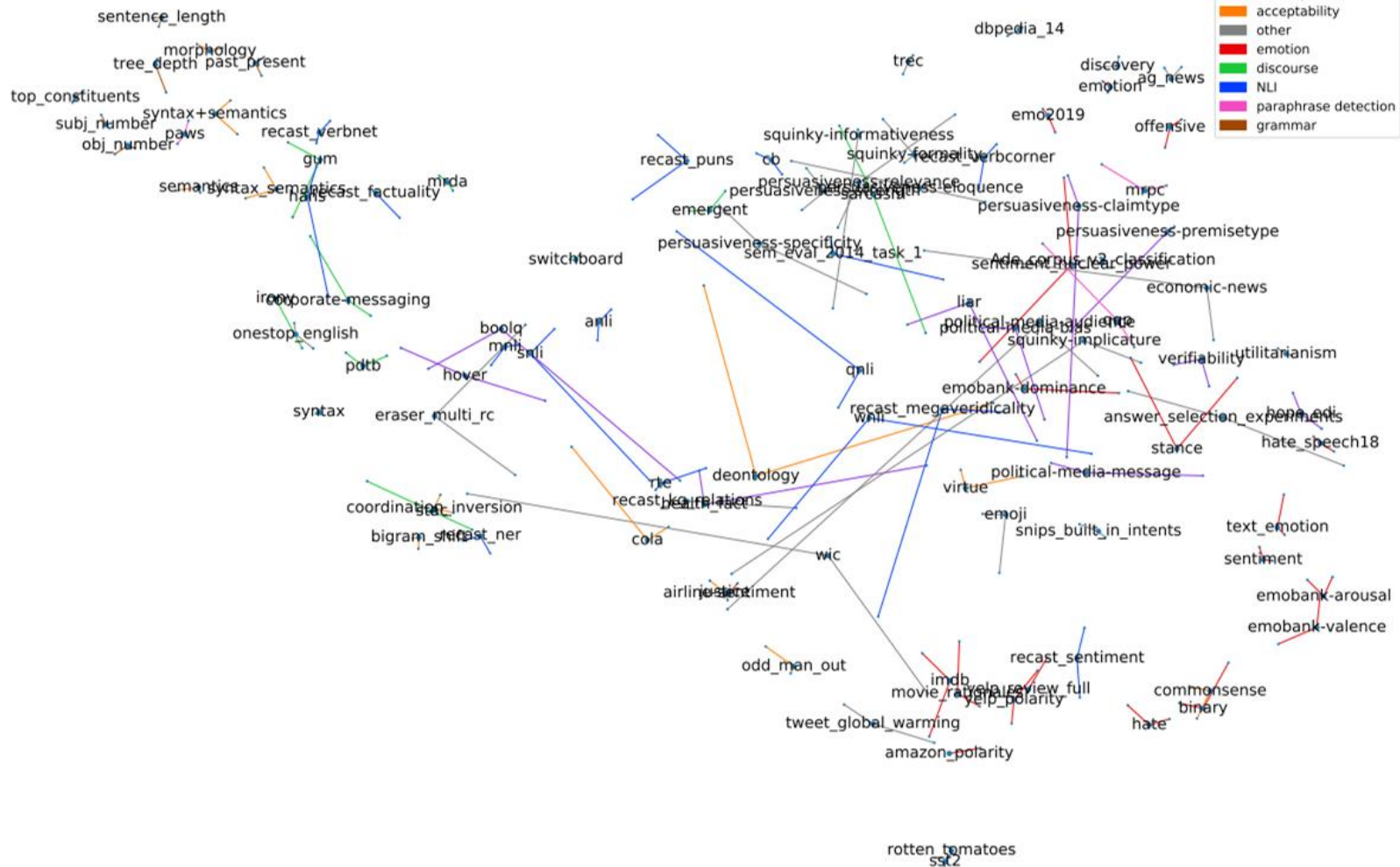
Experimental setting

**Analyses of task embeddings**

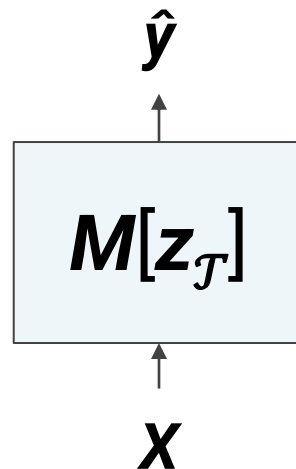
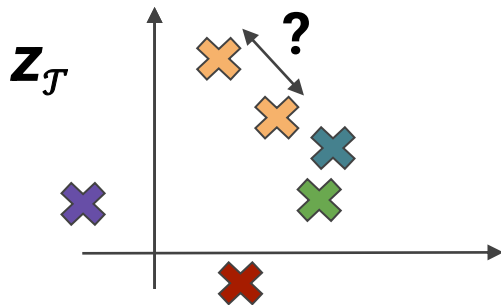
*duplicated task*

2 representations for an identical task  
= proximity ?



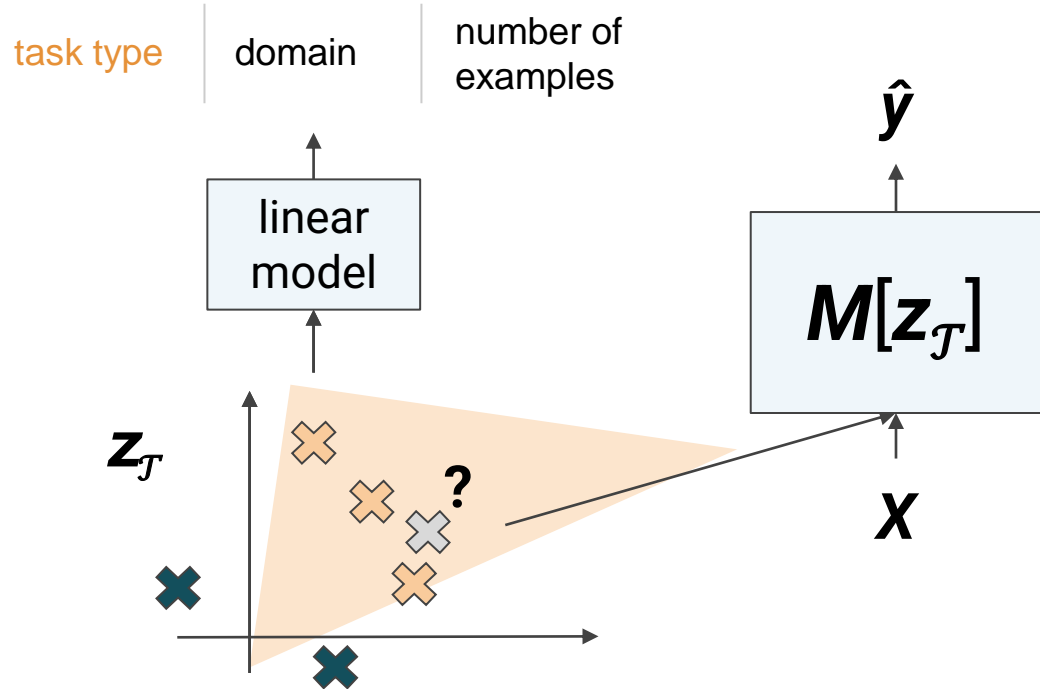


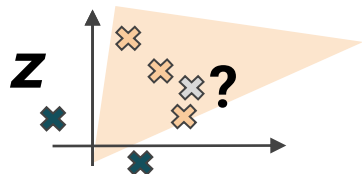
*duplicated task*  
 2 representations for an identical task  
 = proximity ?



Task type	Position stability
Grammar	62.0 ± 3.9
Acceptability	57.1 ± 0.0
Emotion	47.6 ± 2.2
Discourse	45.7 ± 0.0
NLI	37.5 ± 1.0
Other	34.8 ± 0.7
Paraphrase detection	31.5 ± 13.1
Facticity	30.0 ± 4.7
Random embedding	1.0 ± 0.5

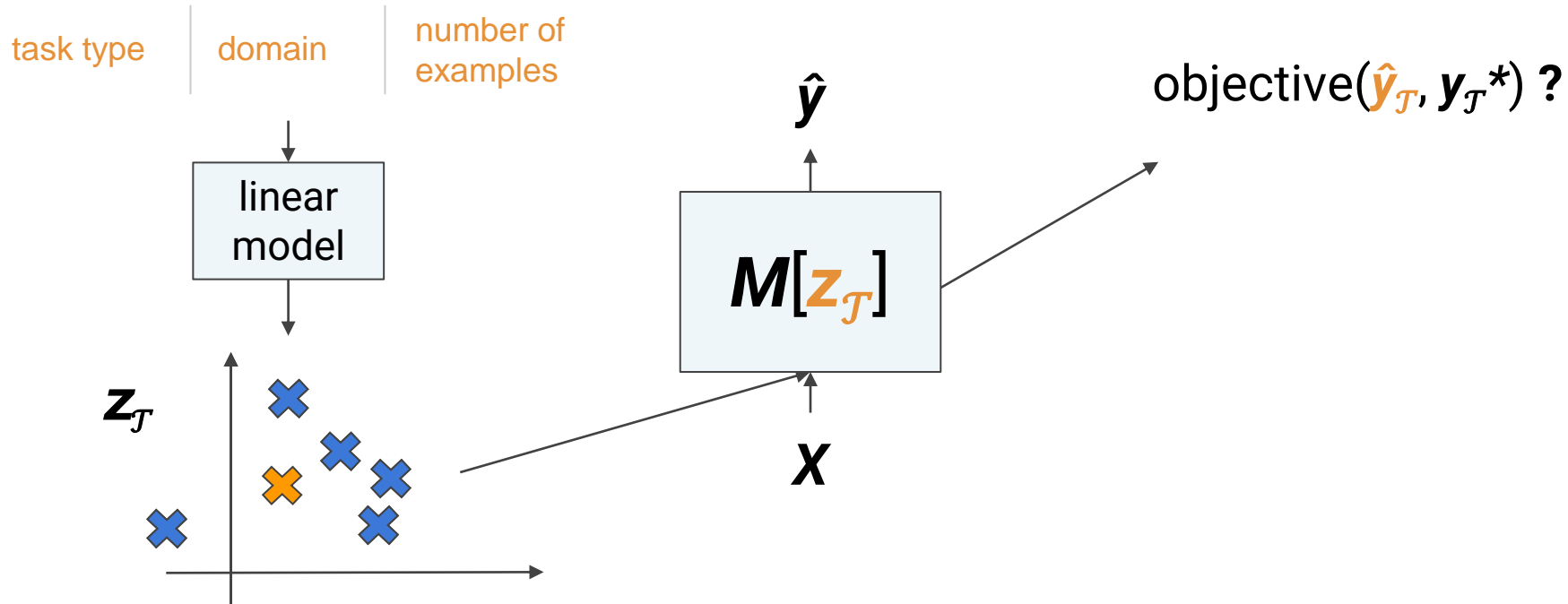




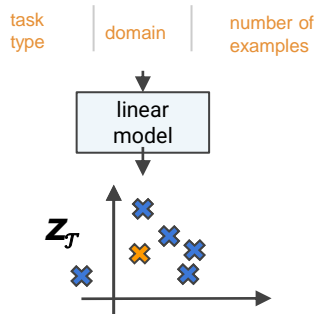


Feature	Domain-Cluster	Num-Rows	Num-Text-Fields	Task-Type	Text-Length
n.a. (Majority Class)	27.8	<b>62.4</b>	63.4	20.8	24.8
Fisher Embedding	37.7	61.3	62.2	35.7	25.6
Learned Features	41.8	51.5	62.2	45.7	32.8
Learned Features $\oplus$ TextEmb	71.5	60.4	76.3	<b>60.4</b>	53.6
TextEmb	<b>78.3</b>	59.3	<b>82.3</b>	50.6	<b>60.3</b>

# Contributions Prediction of task embeddings



# Contributions Prediction of task embeddings



Feature	CoLa	SST2	MRPC	QQP	MNLI	QNLI	RTE	AVG
Single-Task Full-Fine-Tuning (Supervised)	79.2	93.1	75.5	84.7	80.9	88.9	47.3	78.5
Same Task-Type Full Fine-Tuning	73.5	<b>93.6</b>	68.8	55.3	<b>72.7</b>	51.5	<b>70.2</b>	69.4
Features-Aware Task Embeddings - Aspects (ours)	75.4	90.0	<b>70.4</b>	<b>71.1</b>	66.2	56.2	63.7	<b>70.4</b>

# Conclusion

**MetaEval** : 101 aligned text classification datasets

---

Task embeddings lie in **regions** of space

---

**Task embedding** can predict **aspects**

**Aspects** can predict **task embedding**

[damien.sileo@kuleuven.be](mailto:damien.sileo@kuleuven.be)