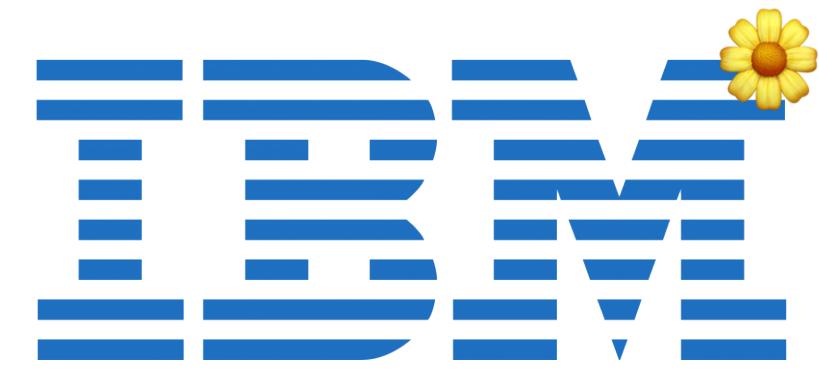


# Hierarchical Pre-training for Sequence Labelling in SpokenDialog

Oral Presentation at WACAI 2021

Emile Chapuis<sup>\*\*</sup>, Pierre Colombo<sup>\*\*</sup>  
Matteo Manica<sup>\*</sup>, Matthieu Labeau<sup>\*</sup>, Chloé Clavel

<sup>\*</sup>Equal contribution



# Hierarchical Pre-training for Sequence Labelling in SpokenDialog

Oral Presentation at WACAI 2021

Emile Chapuis\*, Pierre Colombo\*,  
Matteo Manica, Matthieu Labeau , Chloé Clavel

\*Equal contribution

# Sequence Labelling

---

# Sequence Labelling



**What'd you do, Prison Mike ?**

**I stole. And I robbed.**



**Euhhhh !**



**That is quite the rap sheet, Prison Mike.**



**And I never got caught neither!**



**Well, you are in prison...**

# Sequence Labelling

Dialog/Speech Act



**What'd you do, Prison Mike ?**

**I stole. And I robbed.**



**Euhhhh !**



**That is quite the rap sheet, Prison Mike.**



**And I never got caught neither!**



**Well, you are in prison...**

Question

Answer

Back-channel

Statement opinion

Statement

Appreciation

# Importance of Sequence labelling

# Importance of Sequence labelling

**Types of label considered:**

**Dialog Act**

DIT++ / DAMSL / DiAML

**Emotion and Sentiment**

Polarity (+/-)

6 emotions (+neutrals)

# **Importance of Sequence labelling**

---

**Types of label considered:**

**Dialog Act**

DIT++ / DAMSL / DiAML

**Emotion and Sentiment**

Polarity (+/-)

6 emotions (+neutrals)

**Why are sequence labelling tasks useful?**

**Speaker modelling**

**Dialog State Tracking**

**Avoid generic response problem**

# Importance of Sequence labelling

Generic response problem

Types of label considered:

Dialog Act

DIT++ / DAMSL / DiAML



Emotion and Sentiment

Polarity (+/-)

6 emotions (+neutrals)



Why are sequence labelling tasks useful?

Speaker modelling

Dialog State Tracking

Avoid generic response problem



What'd you do, Prison Mike ?

I don't know



Didnt' you rob someone?

I don't know

## **Limitation of current systems**

---

## **Limitation of current systems**

---

**Current state-of-the art considers deep learning models**

## **Limitation of current systems**

---

**Current state-of-the art considers deep learning models**

**Recurrent encoder + MLP**

## **Limitation of current systems**

---

**Current state-of-the art considers deep learning models**

**Recurrent encoder + MLP**

**Contextual recurrent encoder + MLP**

## **Limitation of current systems**

---

**Current state-of-the art considers deep learning models**

**Recurrent encoder + MLP**

**Contextual recurrent encoder + MLP**

**Contextual recurrent encoder + CRF/Recurrent decoder**

## **Limitation of current systems**

---

**Current state-of-the art considers deep learning models**

**Recurrent encoder + MLP**

**Contextual recurrent encoder + MLP**

**Contextual recurrent encoder + CRF/Recurrent decoder**

**Middle/High size corpora**

**Switchboard Dialog Act (100k + utterances)**

**MRDA (110k + utterances)**

**Dialy Dialog Act (100k + utterances)**

## **Limitation of current systems**

---

**Current state-of-the art considers deep learning models**

**Recurrent encoder + MLP**

**Contextual recurrent encoder + MLP**

**Contextual recurrent encoder + CRF/Recurrent decoder**

**Middle/High size corpora**

**Switchboard Dialog Act (100k + utterances)**

**MRDA (110k + utterances)**

**Dialy Dialog Act (100k + utterances)**

**Not practical but deep learning is data hungry!**

# Contributions

---

## Contributions

---

**Research only consider middle/high size corpora**

## Contributions

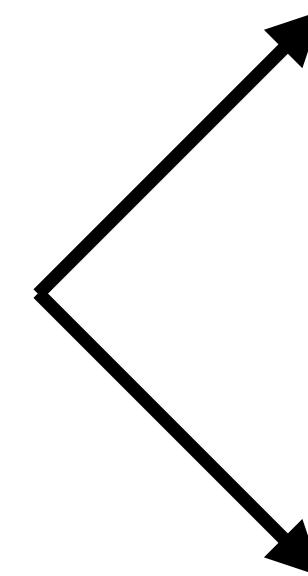
---

Research only consider middle/high size corpora

**SILICONE** (Sequence labellIng evaLuation  
benChmark fOr spoken laNguagE

**Sizes**

**Schema**

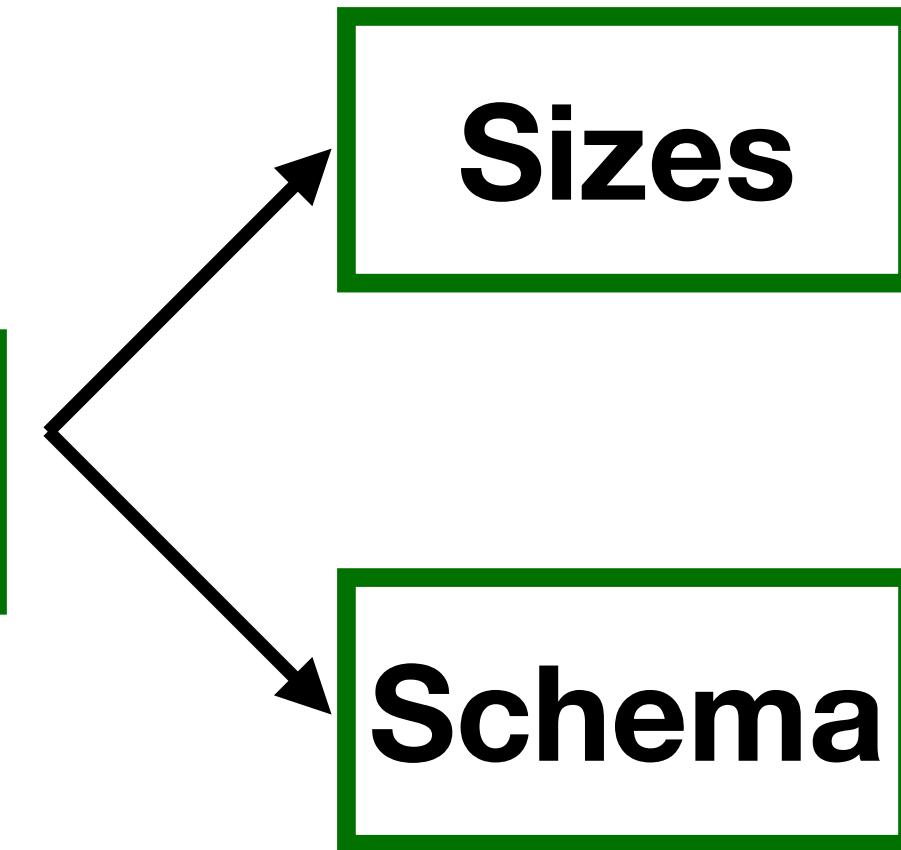


## Contributions

---

Research only consider middle/high size corpora

**SILICONE** (Sequence labellIng evaLuation  
benChmark fOr spoken laNguagE

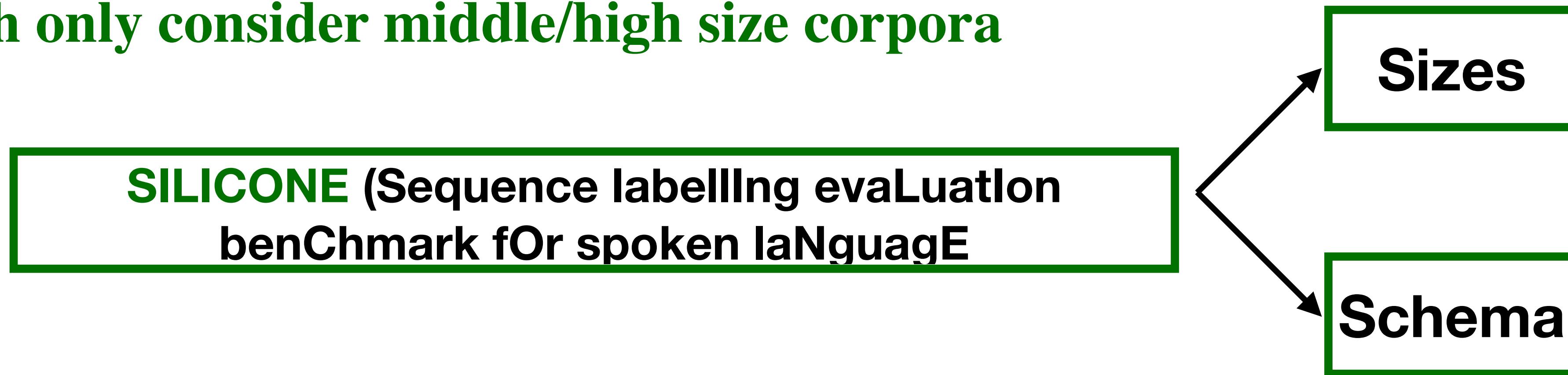


We propose a new model

## Contributions

---

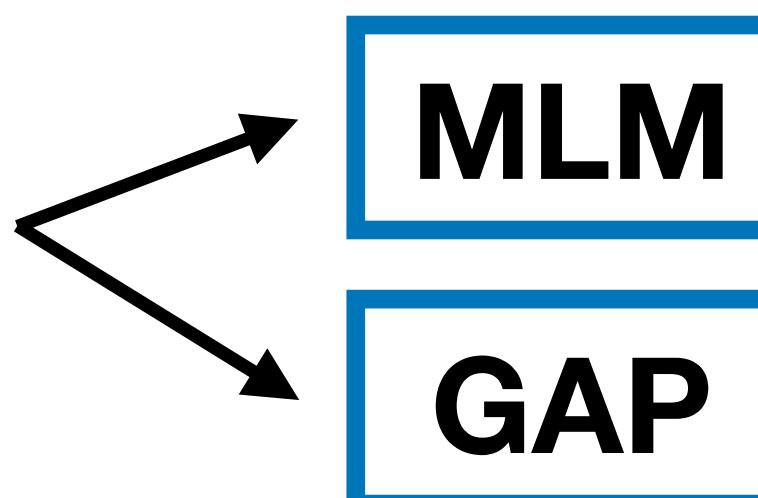
Research only consider middle/high size corpora



We propose a new model

Hierarchical Transformer encoder

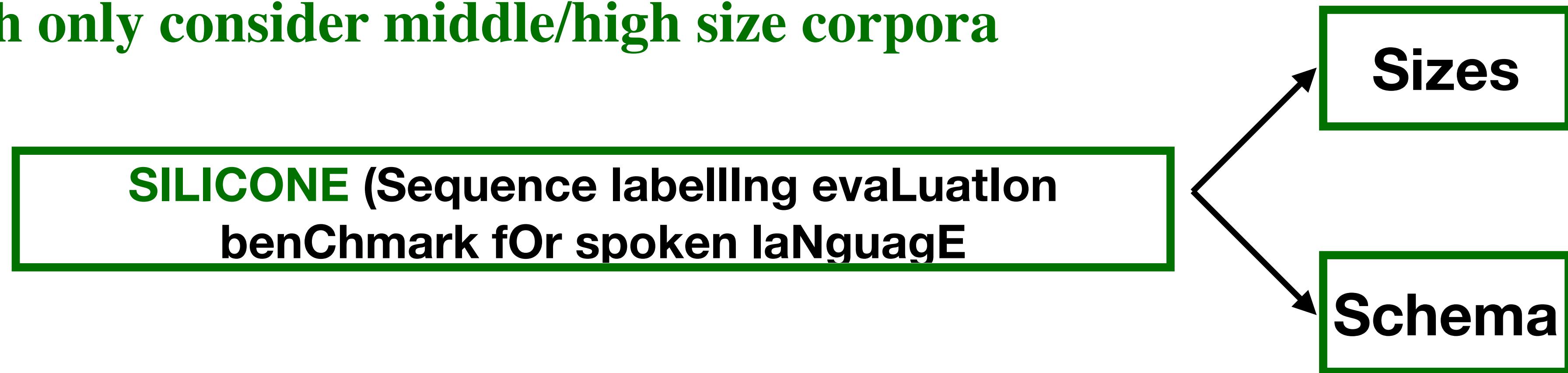
New Pretraining Objectives



## Contributions

---

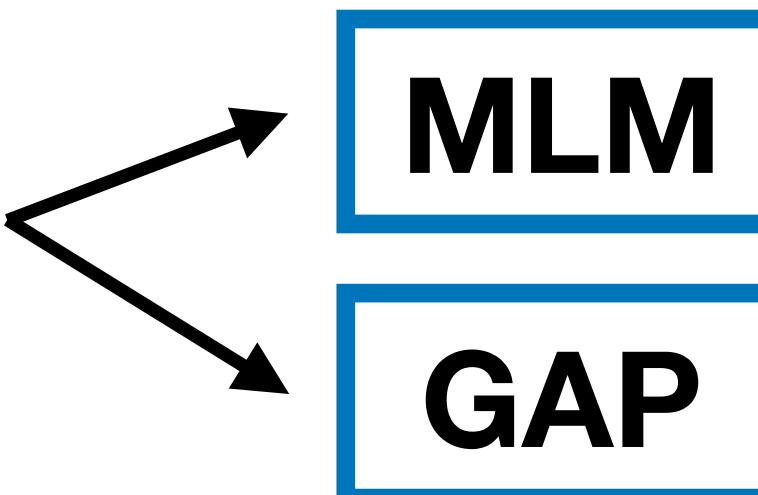
Research only consider middle/high size corpora



We propose a new model

Hierarchical Transformer encoder

New Pretraining Objectives

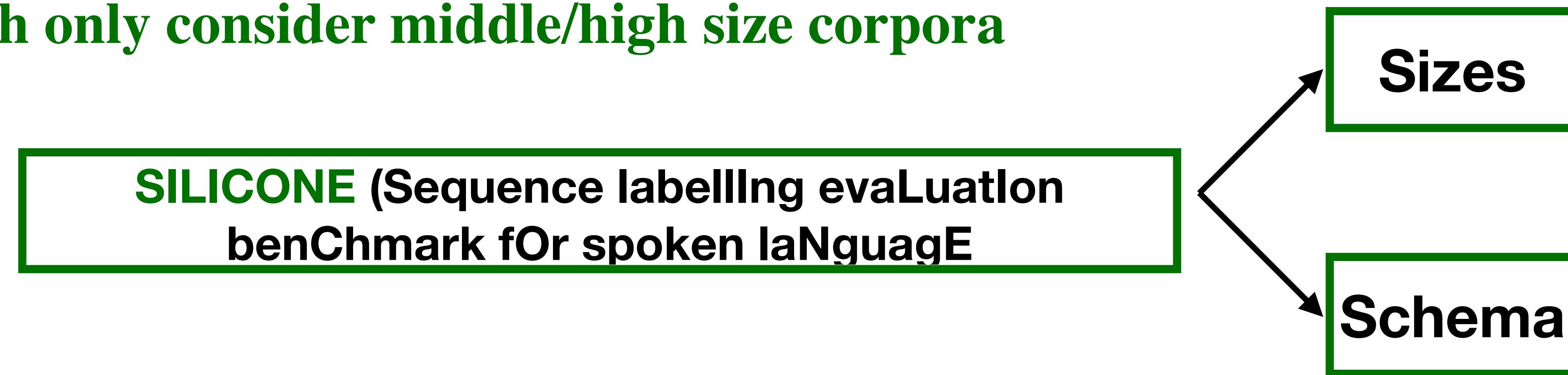


Pretraining Corpora

# Contributions

---

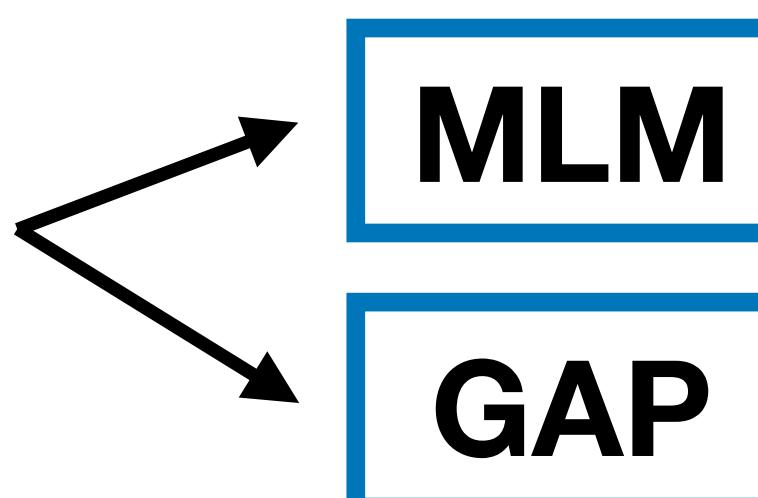
Research only consider middle/high size corpora



We propose a new model

Hierarchical Transformer encoder

New Pretraining Objectives



Pretraining Corpora

OpenSubtitles

Spoken dialogs

Large Scale

Pretraining

# Formalisation

---

# Formalisation

---



**What'd you do, Prison Mike ?**

**I stole. And I robbed.**

**And I kidnapped the president's  
son and held him for ransom.**



**That is quite the rap sheet, Prison Mike.**

**And I never got caught neither!**



**Well, you are in prison...**

# Formalisation

---

$u_1$



**What'd you do, Prison Mike ?**

**I stole. And I robbed.**

**And I kidnapped the president's  
son and held him for ransom.**



$u_2$

$u_3$



**That is quite the rap sheet, Prison Mike.**

**And I never got caught neither!**



$u_4$

$u_5$



**Well, you are in prison...**

# Formalisation

# Conversation $C$

$u_1$



What'd you do, Prison Mike ?

I stole. And I robbed.

And I kidnapped the president's  
son and held him for ransom.



$u_2$

$u_3$



That is quite the rap sheet, Prison Mike.

And I never got caught neither!



$u_4$

$u_5$



Well, you are in prison...

$$C_i = (u_1, u_2, \dots, u_{|C_i|})$$



# Formalisation

---

# Formalisation

---

*u*<sub>1</sub>



Qu'est ce que tu as fait, Prison Mike ?

I stole. And I robbed.



*u*<sub>2</sub>

And I kidnapped the president's  
son and held him for ransom.

*u*<sub>3</sub>



That is quite the rap sheet, Prison Mike.

$\omega_1^4 \quad \omega_2^4 \quad \omega_3^4 \quad \omega_4^4 \quad \omega_5^4 \quad \omega_6^4 \quad \omega_7^4 \quad \omega_8^4 \quad \omega_9^4 \quad \omega_{10}^4$



*u*<sub>4</sub>

Et on ne m'a jamais chopé non plus!

*u*<sub>5</sub>



Well, you are in prison...

# Formalisation

---

$u_1$



Qu'est ce que tu as fait, Prison Mike ?

I stole. And I robbed.



$u_2$

And I kidnapped the president's  
son and held him for ransom.

$u_3$



That is quite the rap sheet, Prison Mike.

$$\omega_1^4 \quad \omega_2^4 \quad \omega_3^4 \quad \omega_4^4 \quad \omega_5^4 \quad \omega_6^4 \quad \omega_7^4 \quad \omega_8^4 \quad \omega_9^4 \quad \omega_{10}^4$$



$u_4$

Et on ne m'a jamais chopé non plus!

$u_5$



Well, you are in prison...

$$u_i^{L_i} = (\omega_1^i, \omega_2^i, \dots, \omega_{|u_i|}^i)$$

# Models

---

# Models

---

Hierarchical Encoder composed of two functions:  $f^u$  and  $f^d$

## Models

---

Hierarchical Encoder composed of two functions:  $f^u$  and  $f^d$

$$\mathcal{E}_{u_i} = f_\theta^u(\omega_1^i, \dots, \omega_{|u_i|}^i),$$

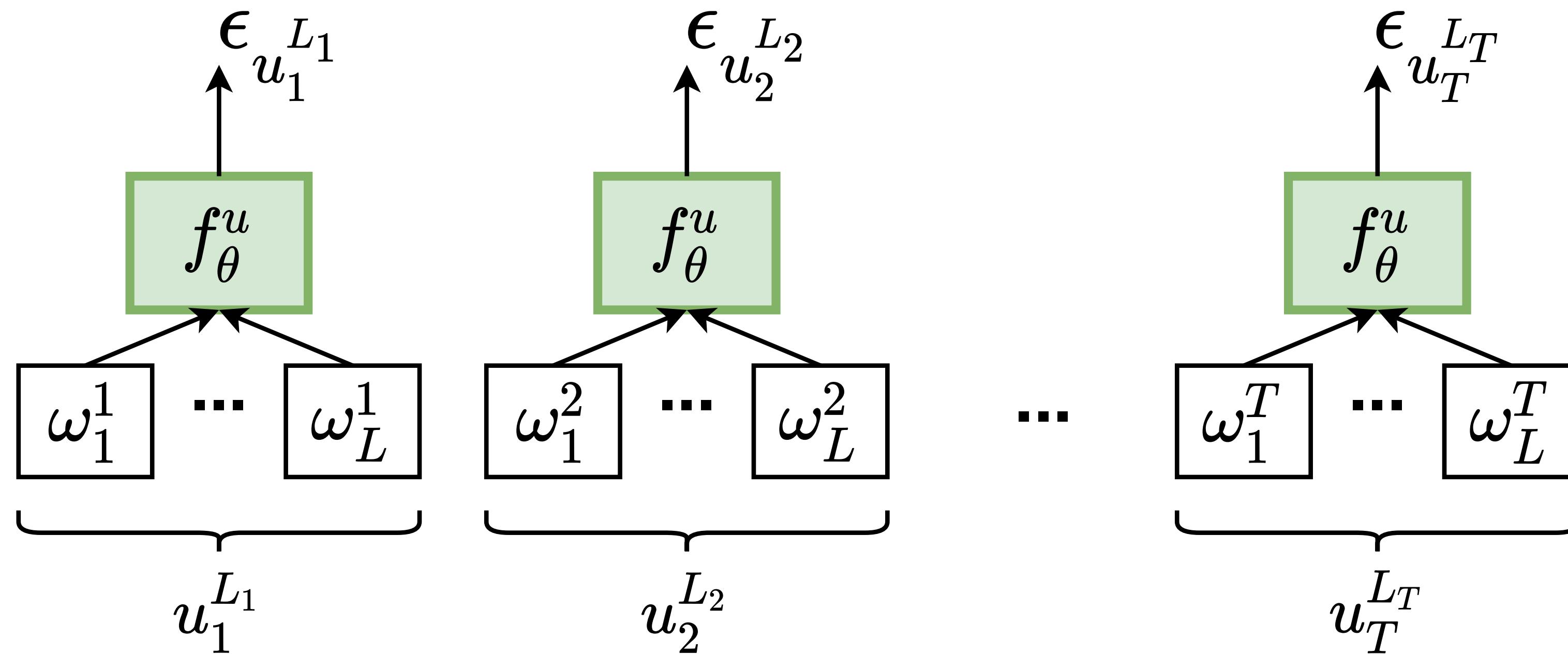
$$\mathcal{E}_{C_j} = f_\theta^d(\mathcal{E}_{u_1}, \dots, \mathcal{E}_{u_{|C_j|}}),$$

# Models

Hierarchical Encoder composed of two functions:  $f^u$  and  $f^d$

$$\mathcal{E}_{u_i^{L_i}} = f_\theta^u(\omega_1^i, \dots, \omega_{|u_i|}^i),$$

$$\mathcal{E}_{C_j} = f_\theta^d(\mathcal{E}_{u_1}, \dots, \mathcal{E}_{u_{|C_j|}}),$$

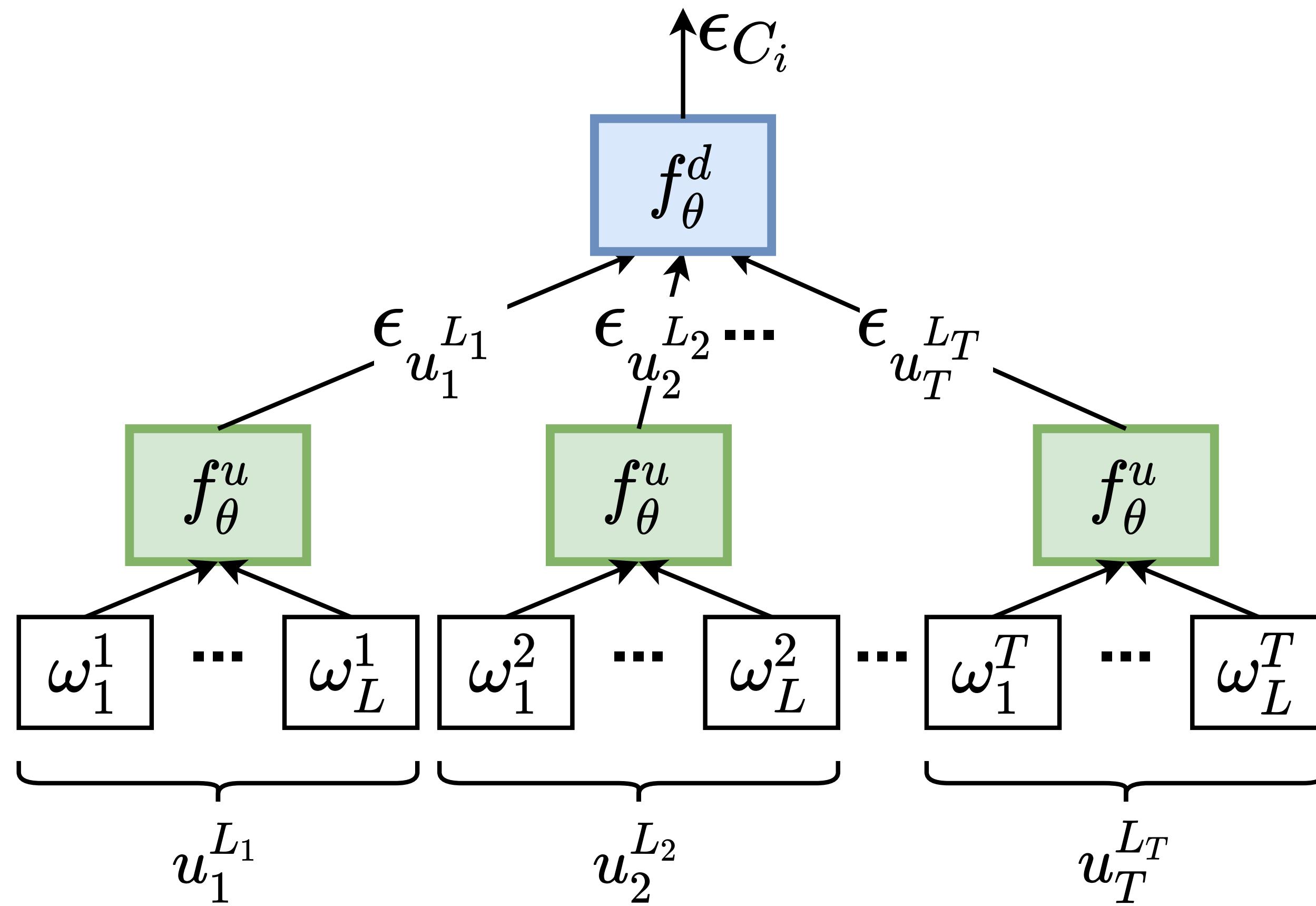


## Models

Hierarchical Encoder composed of two functions:  $f^u$  and  $f^d$

$$\mathcal{E}_{u_i^{L_i}} = f_\theta^u(\omega_1^i, \dots, \omega_{|u_i|}^i),$$

$$\mathcal{E}_{C_j} = f_\theta^d(\mathcal{E}_{u_1}, \dots, \mathcal{E}_{u_{|C_j|}}),$$



## Training Objective

---

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## Models

---

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## Models

---

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

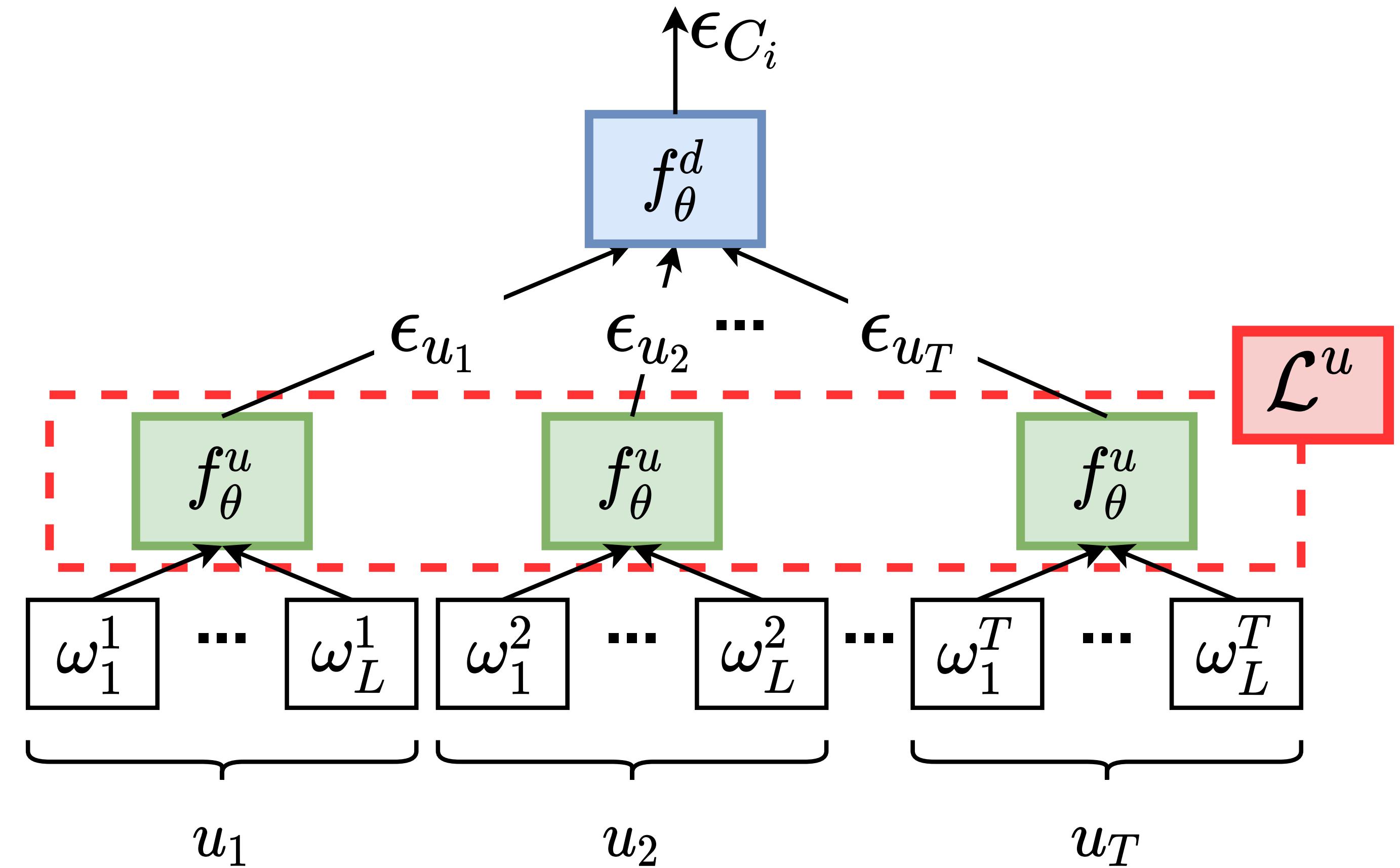
Utterance Level Pretraining

## Models

---

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

### Utterance Level Pretraining



## Losses

---

### Utterance Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## **Losses**

---

**Utterance Level Pretraining**

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

**Masked Utterance Modelling  
(MUM)**

## Losses

---

Utterance Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

Masked Utterance Modelling  
(MUM)

$$p(\Omega \mid \tilde{u}_i^{L_i}) = \prod_{t \in \mathcal{M}_\omega} p_\theta(\omega_t^i \mid \tilde{u}_i^{L_i}).$$

## Losses

---

### Utterance Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

### Masked Utterance Modelling (MUM)

$$p(\boxed{\Omega} \mid \tilde{u}_i^{L_i}) = \prod_{t \in \boxed{\mathcal{M}_{\omega}}} p_{\theta}(\omega_t^i \mid \tilde{u}_i^{L_i}).$$

Set of masked tokens                          Set of masked indices

## Losses

---

### Utterance Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

### Masked Utterance Modelling (MUM)

$\mathcal{U}_3$



That|is|quite|the|rap|sheet,|Prison|Mike|.

$\omega_1^4 \quad \omega_2^4 \quad \omega_3^4 \quad \omega_4^4 \quad \omega_5^4 \quad \omega_6^4 \quad \omega_7^4 \quad \omega_8^4 \quad \omega_9^4 \quad \omega_{10}^4$

$$p(\Omega | \tilde{u}_i^{L_i}) = \prod_{t \in \mathcal{M}_{\omega}} p_{\theta}(\omega_t^i | \tilde{u}_i^{L_i}).$$

Set of masked tokens

Set of masked indices

## Losses

## Utterance Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## Masked Utterance Modelling (MUM)

$\mathcal{U}_3$



That|is|[MASK]|the|rap|[MASK],|Prison|Mike|.

$\omega_1^4 \quad \omega_2^4 \quad \omega_3^4 \quad \omega_4^4 \quad \omega_5^4 \quad \omega_6^4 \quad \omega_7^4 \quad \omega_8^4 \quad \omega_9^4 \quad \omega_{10}^4$

$$p(\Omega | \tilde{u}_i^{L_i}) = \prod_{t \in \mathcal{M}_{\omega}} p_{\theta}(\omega_t^i | \tilde{u}_i^{L_i}).$$

Set of masked tokens

Set of masked indices

Predict → quite sheet

## Models

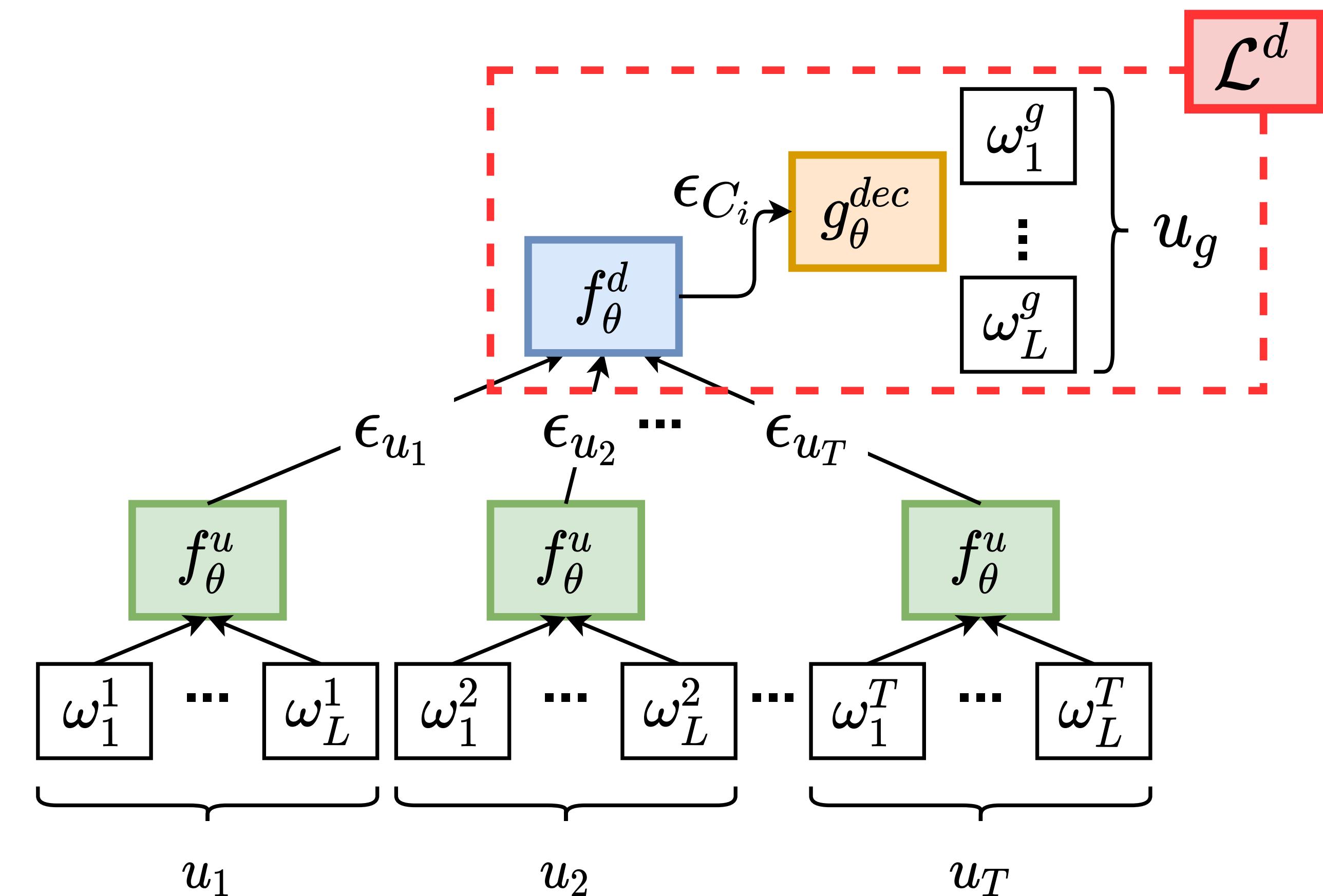
---

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## Models

---

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

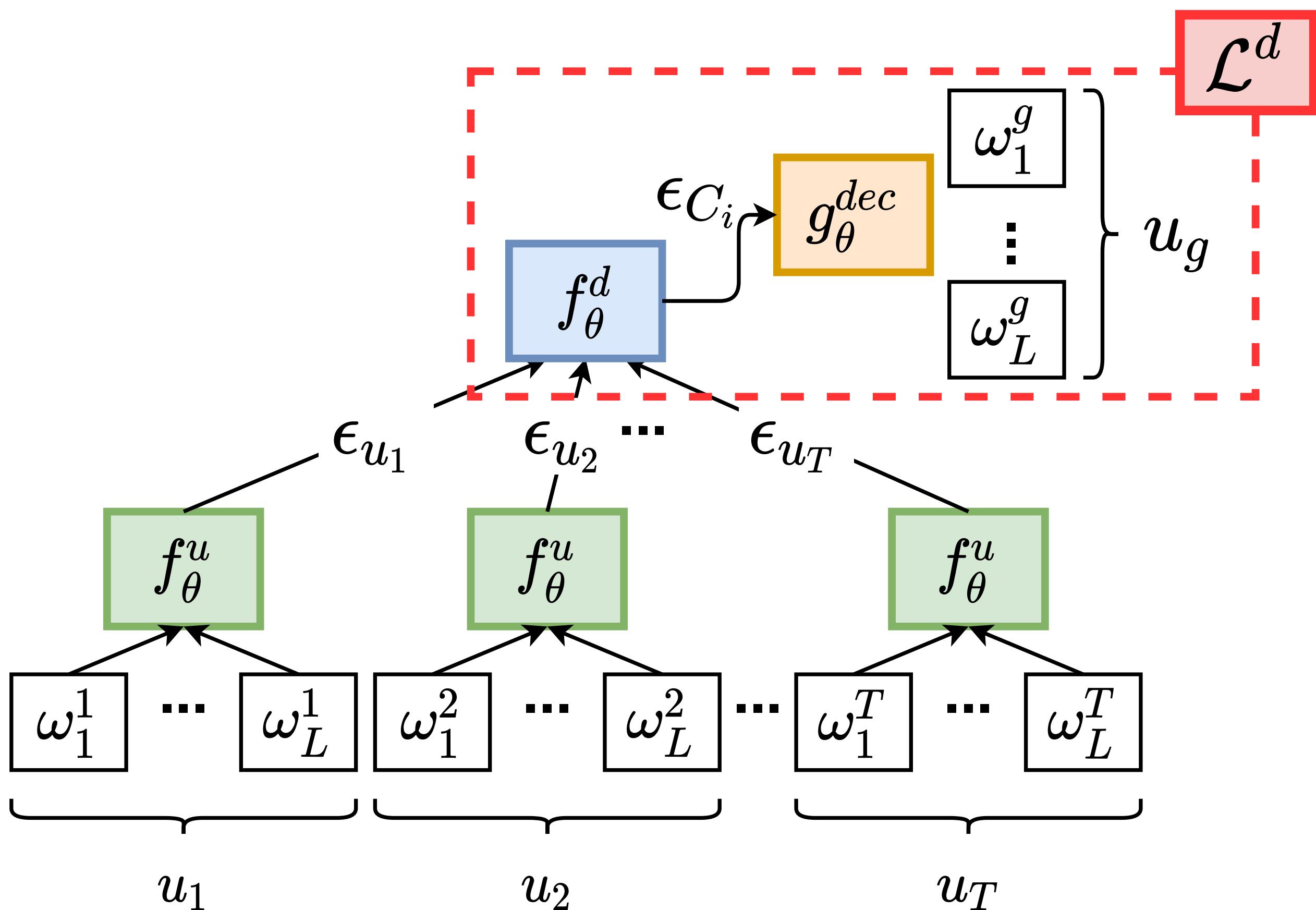


## Models

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

### Dialog Level Pretraining:

- Masked Sequence Generation
- Add a causal decoder



# Models

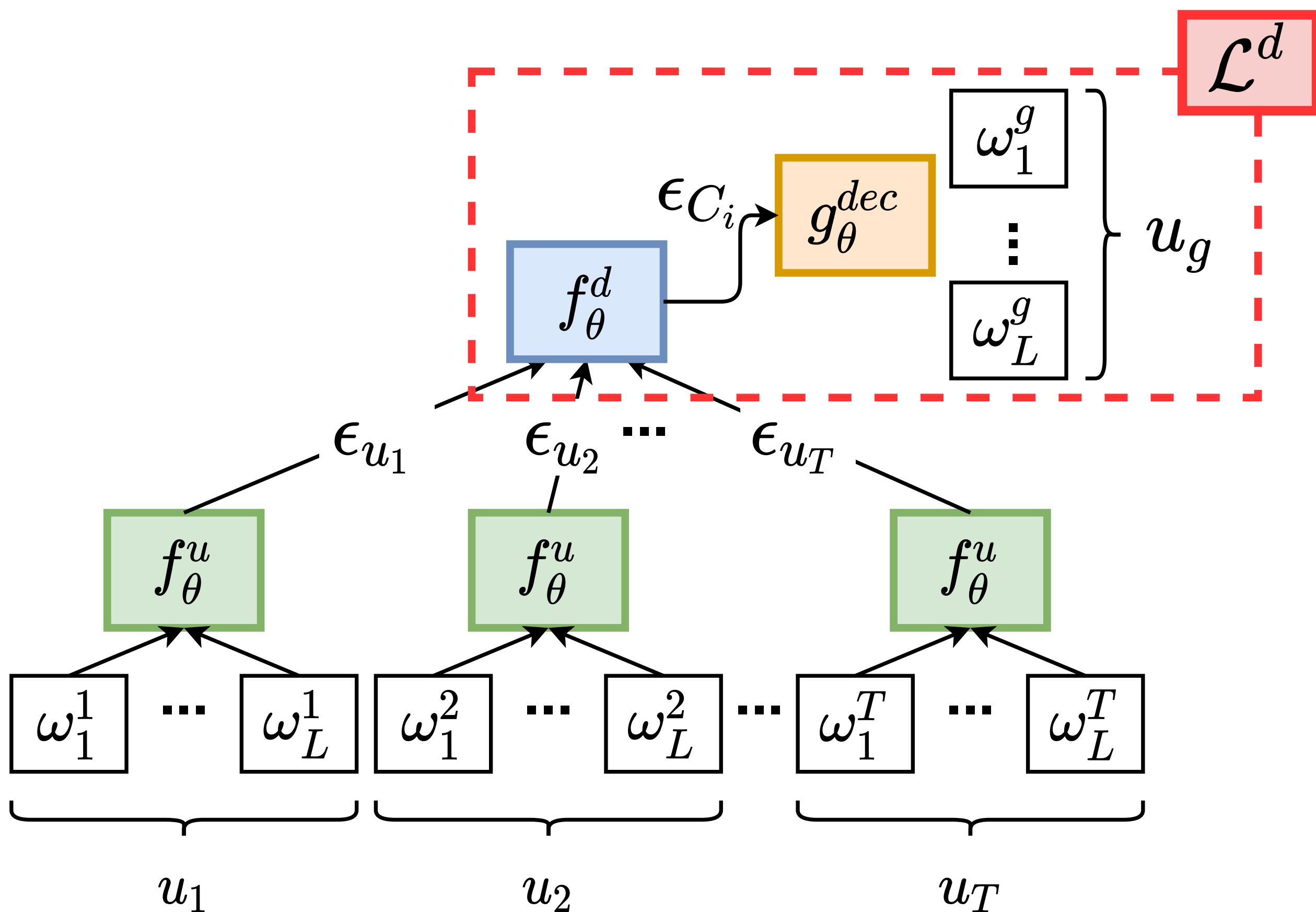
$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## Dialog Level Pretraining:

- Masked Sequence Generation
- Add a causal decoder

## Goal:

- Learnt the **inter-utterance dependancies**
- Train the **2nd level** of the Transformer



## Utterance Level MLM Loss

---

Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## Utterance Level MLM Loss

Masked Language Model Loss  
(MLM)

Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

$$\mathcal{L}_{\text{MLM}}^d(\theta, C_k) = \mathbb{E} \left[ \sum_{j \in m^{C_k}} \sum_{i=1}^{|u_j|} \log(p_\theta(\omega_i^j | \tilde{C}_k)) \right]$$

## Utterance Level MLM Loss

### Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

### Masked Language Model Loss (MLM)

$$\mathcal{L}_{\text{MLM}}^d(\theta, \boxed{C_k}) = \mathbb{E} \left[ \sum_{j \in m^{C_k}} \sum_{i=1}^{|u_j|} \log(p_\theta(\omega_i^j | \tilde{C}_k)) \right]$$

## Utterance Level MLM Loss

### Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

### Masked Language Model Loss (MLM)

$$\mathcal{L}_{\text{MLM}}^d(\theta, \boxed{C_k}) = \mathbb{E} \left[ \sum_{j \in m^{C_k}} \sum_{i=1}^{|u_j|} \log(p_\theta(\omega_i^j | \tilde{C}_k)) \right]$$

## Utterance Level MLM Loss

### Masked Language Model Loss (MLM)

### Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

$$\mathcal{L}_{\text{MLM}}^d(\theta, C_k) = \mathbb{E} \left[ \sum_{j \in m^{C_k}} \sum_{i=1}^{|u_j|} \log(p_\theta(\omega_i^j | \tilde{C}_k)) \right]$$

# Utterance Level MLM Loss

## Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## Masked Language Model Loss (MLM)



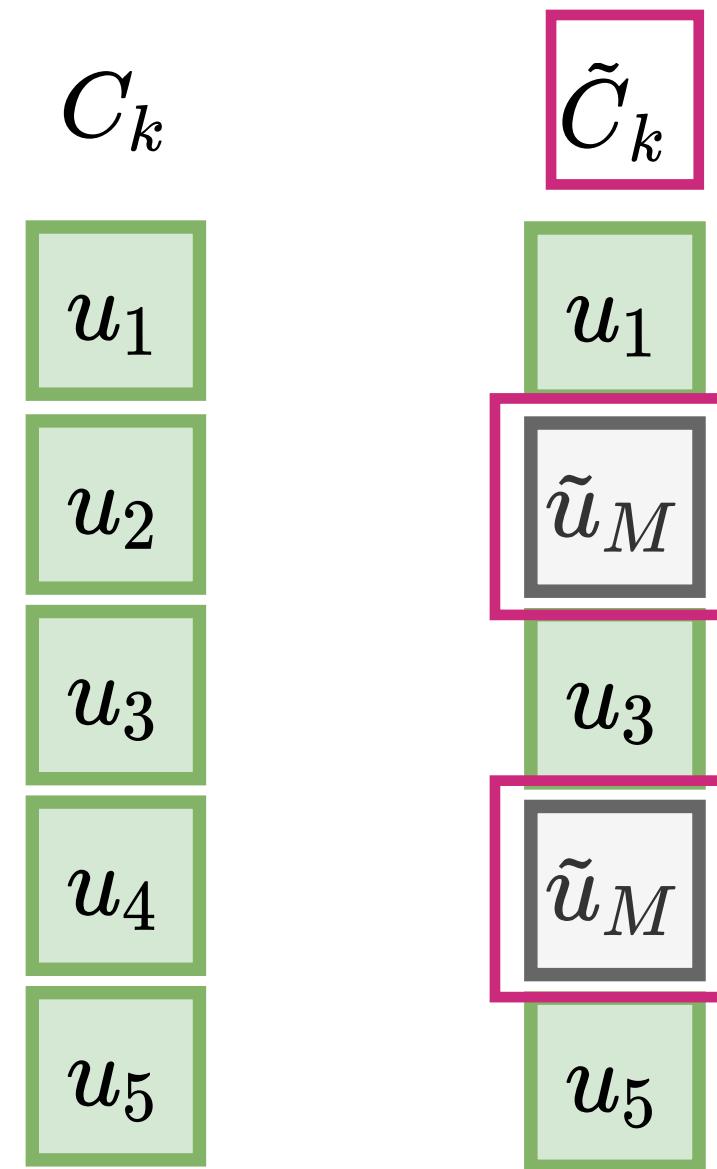
$$\mathcal{L}_{\text{MLM}}^d(\theta, C_k) = \mathbb{E} \left[ \sum_{j \in m^{C_k}} \sum_{i=1}^{|u_j|} \log(p_\theta(\omega_i^j | \tilde{C}_k)) \right]$$

# Utterance Level MLM Loss

## Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## Masked Language Model Loss (MLM)



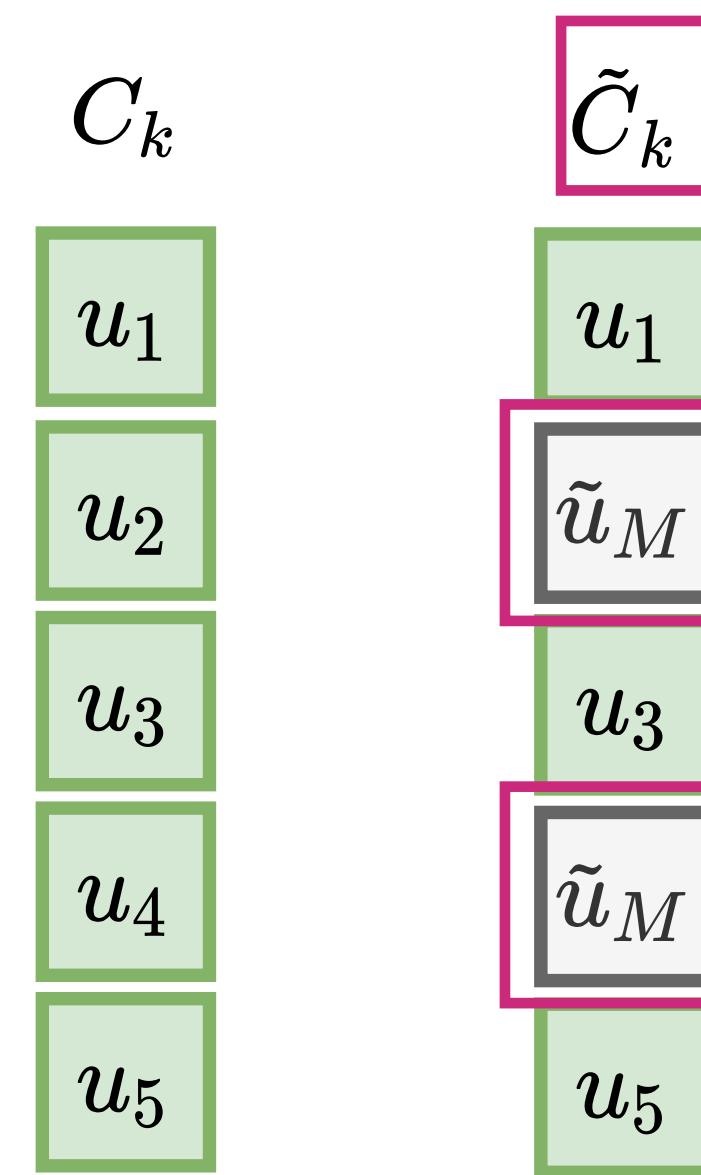
$$\mathcal{L}_{\text{MLM}}^d(\theta, \boxed{C_k}) = \mathbb{E} \left[ \sum_{j \in m^{C_k}} \sum_{i=1}^{|u_j|} \log(p_\theta(\omega_i^j | \boxed{\tilde{C}_k})) \right]$$

# Utterance Level MLM Loss

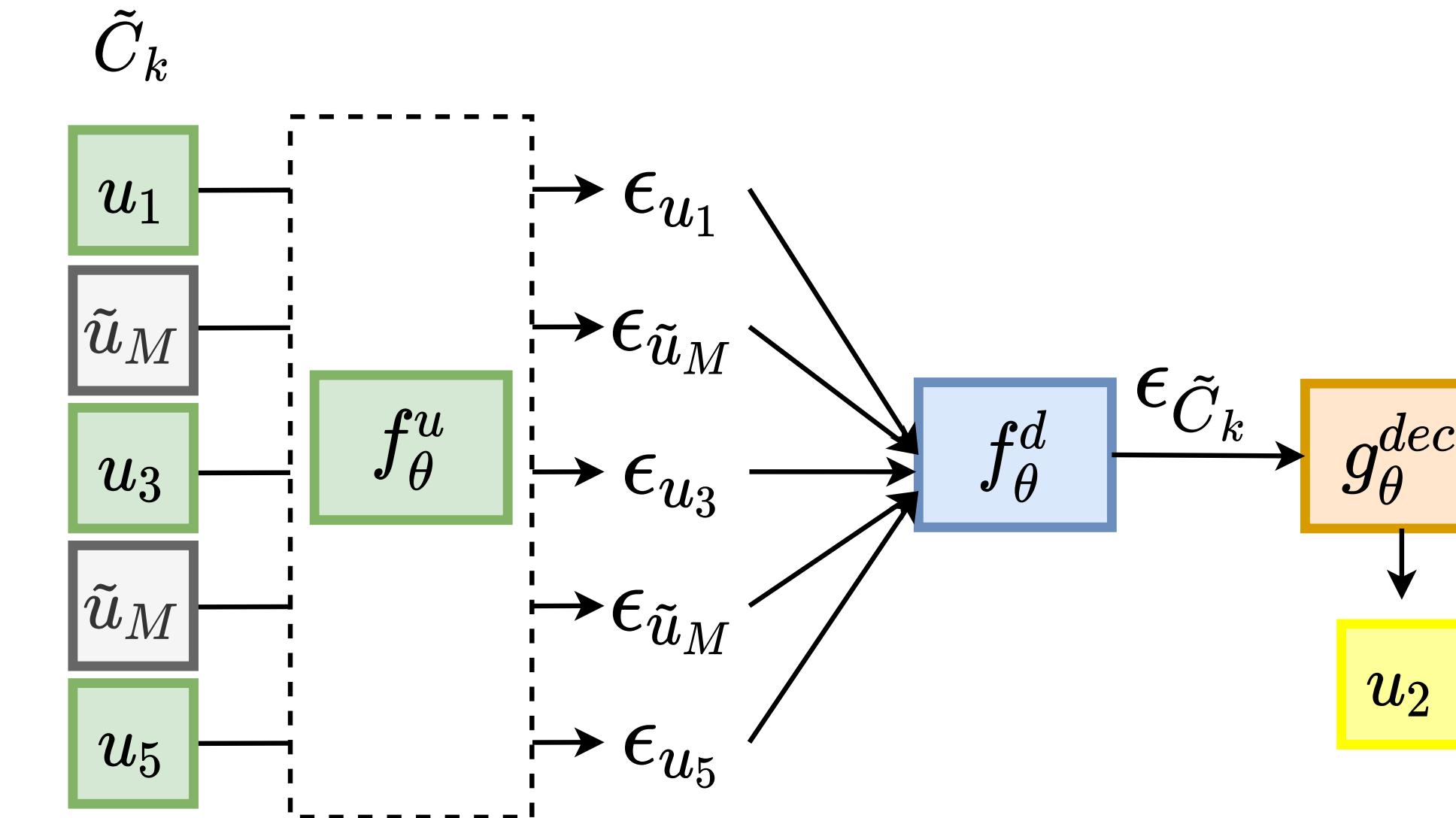
## Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## Masked Language Model Loss (MLM)



$$\mathcal{L}_{\text{MLM}}^d(\theta, \tilde{C}_k) = \mathbb{E} \left[ \sum_{j \in m^{C_k}} \sum_{i=1}^{|u_j|} \log(p_\theta(\omega_i^j | \tilde{C}_k)) \right]$$

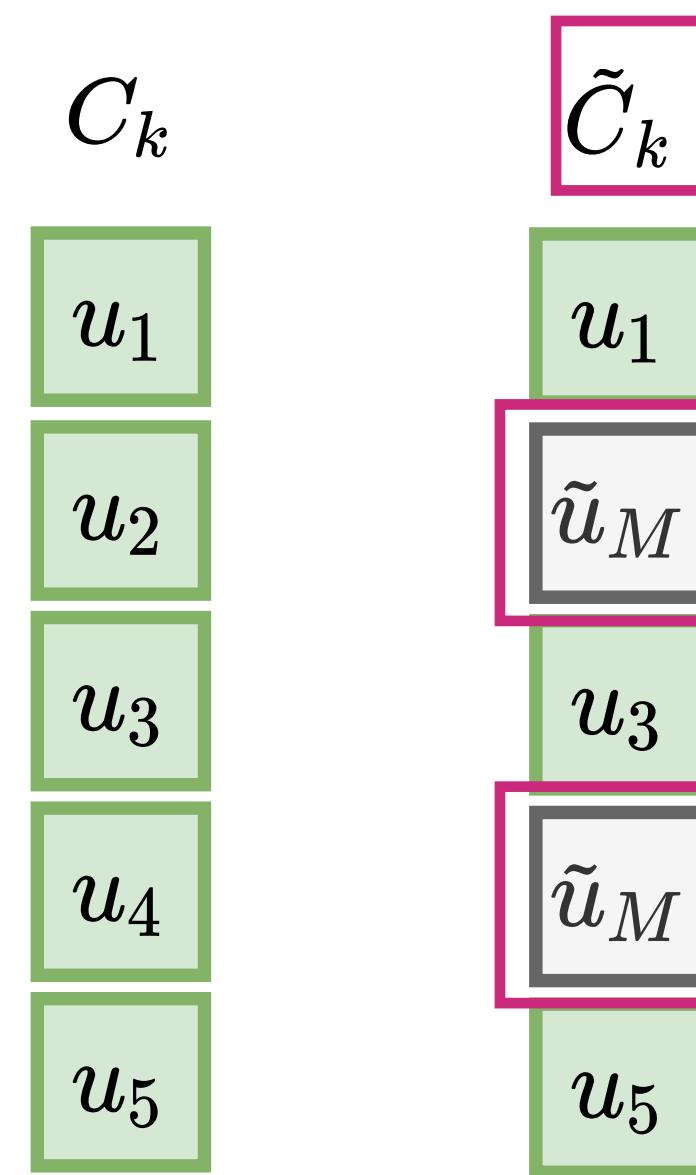


# Utterance Level MLM Loss

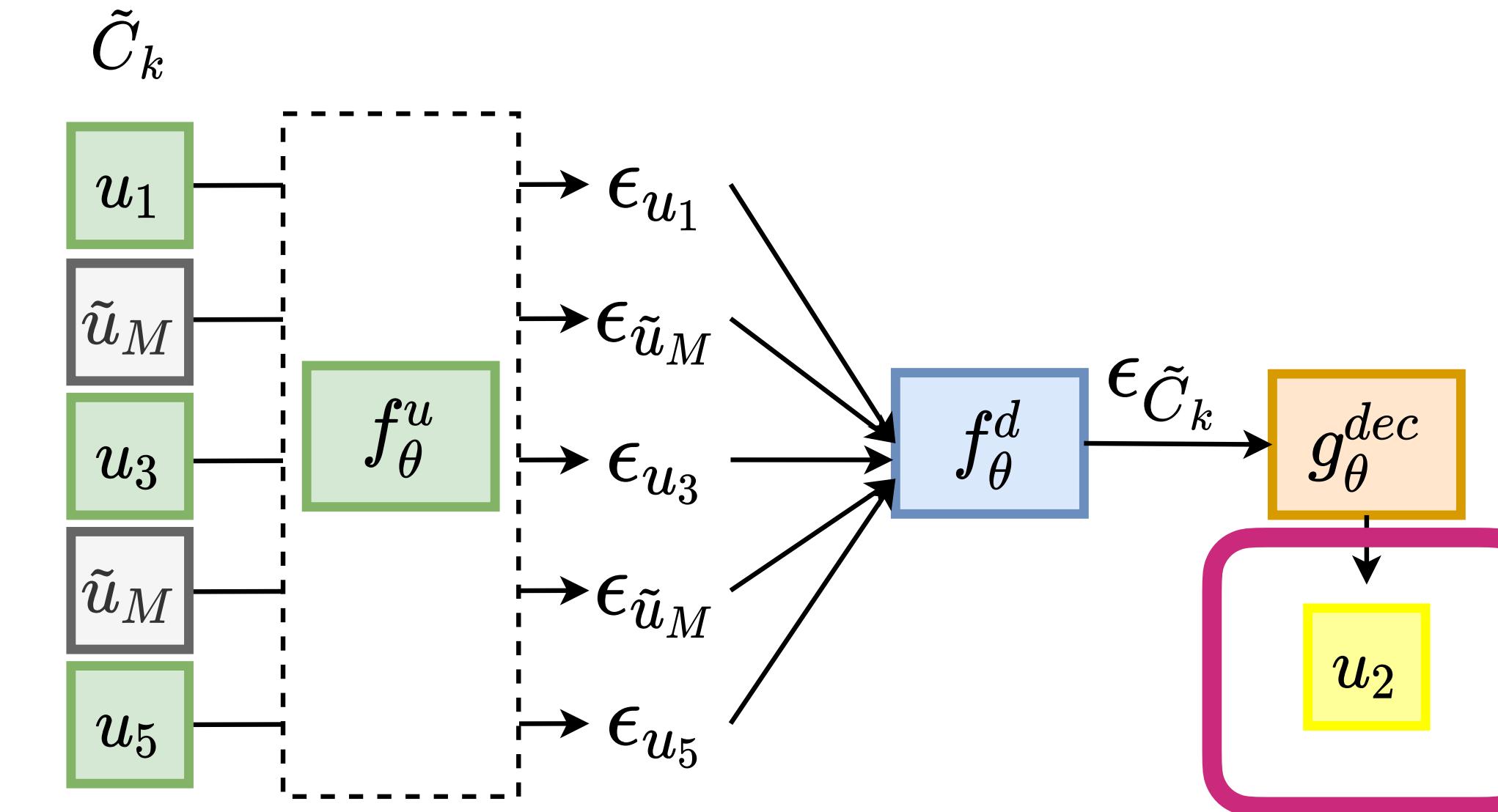
## Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## Masked Language Model Loss (MLM)



$$\mathcal{L}_{\text{MLM}}^d(\theta, \tilde{C}_k) = \mathbb{E} \left[ \sum_{j \in m^{C_k}} \sum_{i=1}^{|u_j|} \log(p_\theta(\omega_i^j | \tilde{C}_k)) \right]$$

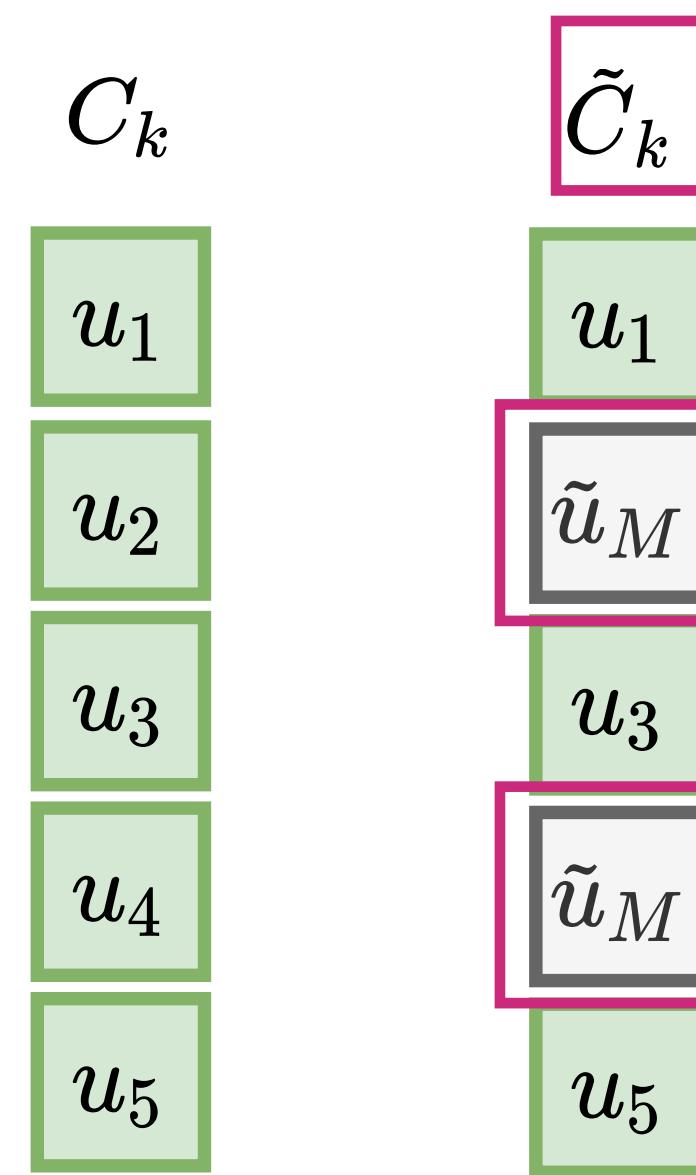


# Utterance Level MLM Loss

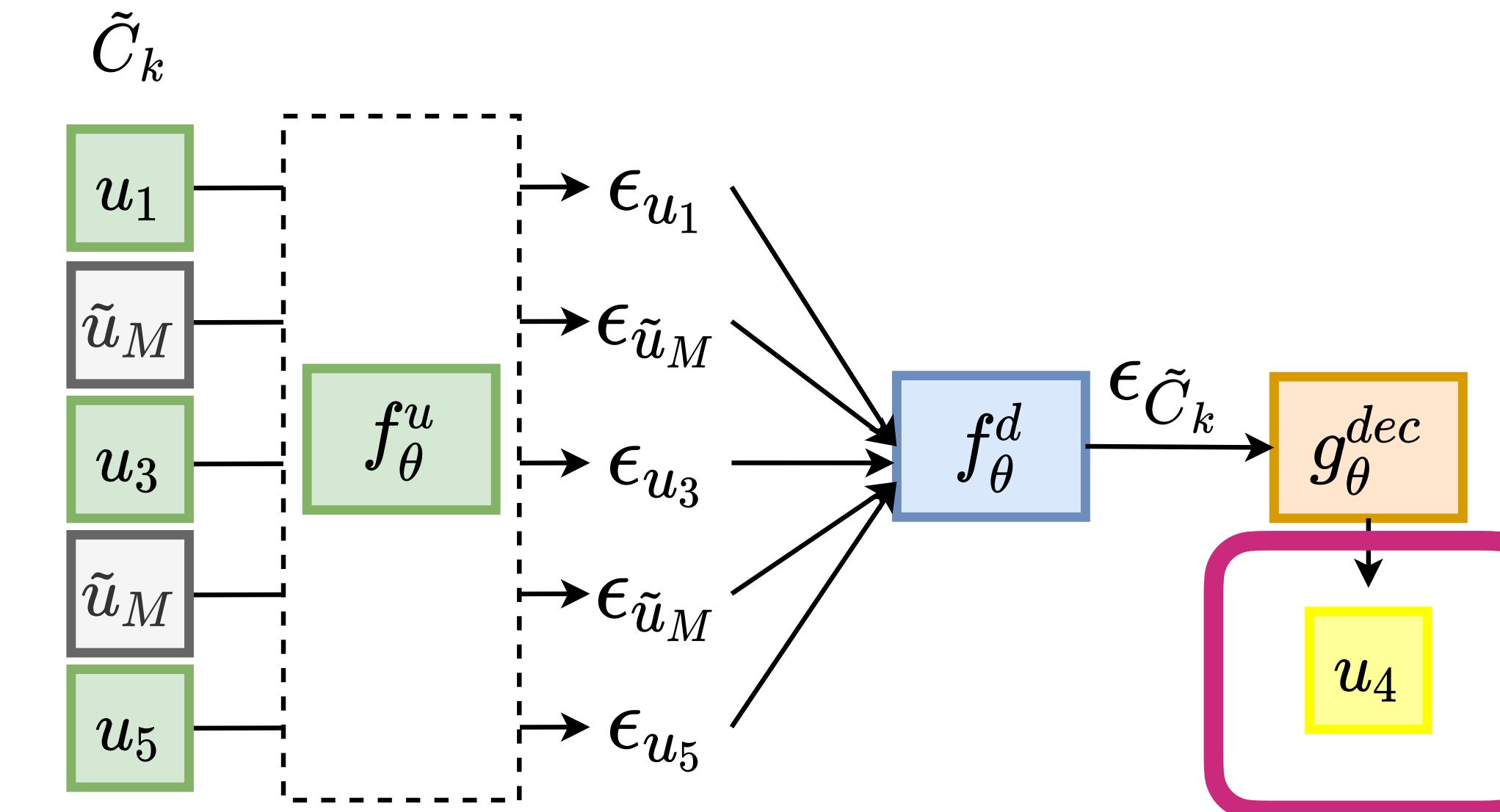
Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

## Masked Language Model Loss (MLM)



$$\mathcal{L}_{\text{MLM}}^d(\theta, \tilde{C}_k) = \mathbb{E} \left[ \sum_{j \in m^{C_k}} \sum_{i=1}^{|u_j|} \log(p_\theta(\omega_i^j | \tilde{C}_k)) \right]$$



# SILICONE

---

**Research only consider middle/high size corpora**

# SILICONE

---

**Research only consider middle/high size corpora**

Corpus	<i>Train</i>	<i>Val</i>	<i>Test</i>	Utt.	<i>Labels</i>	Task	Utt./ <i>Labels</i>
SwDA*	1k	100	11	200k	42	DA	4.8k
MRDA*	56	6	12	110k	5	DA	2.6k
DyDA <sub>a</sub>	11k	1k	1k	102k	4	DA	25.5k
MT*	121	22	25	36k	12	DA	3k
Oasis*	508	64	64	15k	42	DA	357
DyDA <sub>e</sub>	11k	1k	1k	102k	7	E	2.2k
MELD <sub>S</sub> *	934	104	280	13k	3	S	4.3k
MELD <sub>E</sub> *	934	104	280	13k	7	S	1.8k
IEMO	108	12	31	10k	6	E	1.7k
SEM	62	7	10	5,6k	3	S	1.9k

# SILICONE

Research only consider middle/high size corpora

Dialog Acts

Corpus	Train	Val	Test	Utt.	Labels	Task	Utt./Labels
SwDA*	1k	100	11	200k	42	DA	4.8k
MRDA*	56	6	12	110k	5	DA	2.6k
DyDA <sub>a</sub>	11k	1k	1k	102k	4	DA	25.5k
MT*	121	22	25	36k	12	DA	3k
Oasis*	508	64	64	15k	42	DA	357
DyDA <sub>e</sub>	11k	1k	1k	102k	7	E	2.2k
MELD <sub>S</sub> *	934	104	280	13k	3	S	4.3k
MELD <sub>E</sub> *	934	104	280	13k	7	S	1.8k
IEMO	108	12	31	10k	6	E	1.7k
SEM	62	7	10	5,6k	3	S	1.9k

# SILICONE

Research only consider middle/high size corpora

Dialog Acts



Corpus	Train	Val	Test	Utt.	Labels	Task	Utt./Labels
SwDA*	1k	100	11	200k	42	DA	4.8k
MRDA*	56	6	12	110k	5	DA	2.6k
DyDA <sub>a</sub>	11k	1k	1k	102k	4	DA	25.5k
MT*	121	22	25	36k	12	DA	3k
Oasis*	508	64	64	15k	42	DA	357
<hr/>							
DyDA <sub>e</sub>	11k	1k	1k	102k	7	E	2.2k
MELD <sub>S</sub> *	934	104	280	13k	3	S	4.3k
MELD <sub>e</sub> *	934	104	280	13k	7	S	1.8k
IEMO	108	12	31	10k	6	E	1.7k
SEM	62	7	10	5,6k	3	S	1.9k

Emotions  
&  
Sentiments



## **General Performances**

---

### **Results on SILICONE**

# General Performances

---

## Results on SILICONE

	Avg	SwDA	MRDA	DyDADA	MT	Oasis	DyDA <sub>e</sub>	MELD <sub>s</sub>	MELD <sub>e</sub>	IEMO	SEM
BERT-4layers	70.4	77.8	90.7	79.0	88.4	66.8	90.3	55.3	53.4	43.0	58.8
BERT	72.8	79.2	90.7	<b>82.6</b>	88.2	66.9	91.9	59.3	<b>61.4</b>	<b>45.0</b>	62.7
$\mathcal{H}\mathcal{R}$	69.8	77,5	90,9	80,1	82,8	64,3	91.5	59,3	59.9	40.3	51.1
$\mathcal{H}\mathcal{T}(\theta_{MLM}^{u,d})$ (TINY)	73.3	<b>79.3</b>	92.0	80.1	90.0	68,3	92.5	62.6	59.9	42.0	66.6
$\mathcal{H}\mathcal{T}(\theta_{GAP}^d)$ (TINY)	71.6	78.6	91.8	78.1	89.3	64.1	91.6	60.5	55.7	42.2	63.9
$\mathcal{H}\mathcal{T}(\theta_{MLM}^{u,d})$ (SMALL)	<b>74.3</b>	79.2	<b>92.4</b>	81.5	<b>90.6</b>	<b>69.4</b>	<b>92.7</b>	<b>64.1</b>	60.1	<b>45.0</b>	<b>68.2</b>

# General Performances

## Results on SILICONE

	Avg	SwDA	MRDA	DyDADA	MT	Oasis	DyDA <sub>e</sub>	MELD <sub>s</sub>	MELD <sub>e</sub>	IEMO	SEM
BERT-4layers	70.4	77.8	90.7	79.0	88.4	66.8	90.3	55.3	53.4	43.0	58.8
BERT	72.8	79.2	90.7	<b>82.6</b>	88.2	66.9	91.9	59.3	<b>61.4</b>	<b>45.0</b>	62.7
$\mathcal{H}\mathcal{R}$	69.8	77.5	90.9	80.1	82.8	64.3	91.5	59.3	59.9	40.3	51.1
$\mathcal{HT}(\theta_{MLM}^{u,d})$ (TINY)	73.3	<b>79.3</b>	92.0	80.1	90.0	68.3	92.5	62.6	59.9	42.0	66.6
$\mathcal{HT}(\theta_{GAP}^d)$ (TINY)	71.6	78.6	91.8	78.1	89.3	64.1	91.6	60.5	55.7	42.2	63.9
$\mathcal{HT}(\theta_{MLM}^{u,d})$ (SMALL)	<b>74.3</b>	79.2	<b>92.4</b>	81.5	<b>90.6</b>	<b>69.4</b>	<b>92.7</b>	<b>64.1</b>	60.1	<b>45.0</b>	<b>68.2</b>

# General Performances

## Results on SILICONE

	Avg	SwDA	MRDA	DyDADA	MT	Oasis	DyDA <sub>e</sub>	MELD <sub>s</sub>	MELD <sub>e</sub>	IEMO	SEM
BERT-4layers	70.4	77.8	90.7	79.0	88.4	66.8	90.3	55.3	53.4	43.0	58.8
BERT	72.8	79.2	90.7	<b>82.6</b>	88.2	66.9	91.9	59.3	<b>61.4</b>	<b>45.0</b>	62.7
$\mathcal{H}\mathcal{R}$	69.8	77.5	90.9	80.1	82.8	64.3	91.5	59.3	59.9	40.3	51.1
$\mathcal{HT}(\theta_{MLM}^{u,d})$ (TINY)	73.3	<b>79.3</b>	92.0	80.1	90.0	68.3	92.5	62.6	59.9	42.0	66.6
$\mathcal{HT}(\theta_{GAP}^d)$ (TINY)	71.6	78.6	91.8	78.1	89.3	64.1	91.6	60.5	55.7	42.2	63.9
$\mathcal{HT}(\theta_{MLM}^{u,d})$ (SMALL)	74.3	79.2	<b>92.4</b>	81.5	<b>90.6</b>	<b>69.4</b>	<b>92.7</b>	<b>64.1</b>	60.1	<b>45.0</b>	<b>68.2</b>

## General Performances

### Results on SILICONE

	Avg	SwDA	MRDA	DyDADA	MT	Oasis	DyDA <sub>e</sub>	MELD <sub>s</sub>	MELD <sub>e</sub>	IEMO	SEM
BERT-4layers	70.4	77.8	90.7	79.0	88.4	66.8	90.3	55.3	53.4	43.0	58.8
BERT	72.8	79.2	90.7	<b>82.6</b>	88.2	66.9	91.9	59.3	<b>61.4</b>	<b>45.0</b>	62.7
$\mathcal{H}\mathcal{R}$	69.8	77.5	90.9	80.1	82.8	64.3	91.5	59.3	59.9	40.3	51.1
$\mathcal{HT}(\theta_{MLM}^{u,d})$ (TINY)	73.3	<b>79.3</b>	92.0	80.1	90.0	68.3	92.5	62.6	59.9	42.0	66.6
$\mathcal{HT}(\theta_{GAP}^d)$ (TINY)	71.6	78.6	91.8	78.1	89.3	64.1	91.6	60.5	55.7	42.2	63.9
$\mathcal{HT}(\theta_{MLM}^{u,d})$ (SMALL)	74.3	79.2	<b>92.4</b>	81.5	<b>90.6</b>	<b>69.4</b>	<b>92.7</b>	<b>64.1</b>	60.1	<b>45.0</b>	<b>68.2</b>

Increasing model size is a good idea!

## General Performances

### Results on SILICONE

	Avg	SwDA	MRDA	DyDADA	MT	Oasis	DyDA <sub>e</sub>	MELD <sub>s</sub>	MELD <sub>e</sub>	IEMO	SEM
BERT-4layers	70.4	77.8	90.7	79.0	88.4	66.8	90.3	55.3	53.4	43.0	58.8
BERT	72.8	79.2	90.7	<b>82.6</b>	88.2	66.9	91.9	59.3	<b>61.4</b>	<b>45.0</b>	62.7
$\mathcal{H}\mathcal{R}$	69.8	77.5	90.9	80.1	82.8	64.3	91.5	59.3	59.9	40.3	51.1
$\mathcal{HT}(\theta_{MLM}^{u,d})$ (TINY)	73.3	<b>79.3</b>	92.0	80.1	90.0	68.3	92.5	62.6	59.9	42.0	66.6
$\mathcal{HT}(\theta_{GAP}^d)$ (TINY)	71.6	78.6	91.8	78.1	89.3	64.1	91.6	60.5	55.7	42.2	63.9
$\mathcal{HT}(\theta_{MLM}^{u,d})$ (SMALL)	74.3	79.2	<b>92.4</b>	81.5	<b>90.6</b>	<b>69.4</b>	<b>92.7</b>	<b>64.1</b>	60.1	<b>45.0</b>	<b>68.2</b>

Our hierarchical pertaining helps!

Increasing model size is a good idea!

# Results I

---

## Results I

---

### Decoder Choice

## Results I

---

### Decoder Choice

	Avg	Avg DA	Avg E/S
BERT (+MLP)	72,8	81.5	64.0
BERT (+GRU)	69.9	80.4	59.3
BERT (+CRF)	72.8	81.5	64.1
$\mathcal{H}\mathcal{R}$ (+MLP)	69.8	79.1	60.4
$\mathcal{H}\mathcal{R}$ (+GRU)	67.6	79.4	55.7
$\mathcal{H}\mathcal{R}$ (+CRF)	70.5	80.3	60.7

## Results I

---

### Decoder Choice

	Avg	Avg DA	Avg E/S
BERT (+MLP)	72,8	81.5	64.0
BERT (+GRU)	69.9	80.4	59.3
BERT (+CRF)	72.8	81.5	64.1
$\mathcal{H}\mathcal{R}$ (+MLP)	69.8	79.1	60.4
$\mathcal{H}\mathcal{R}$ (+GRU)	67.6	79.4	55.7
$\mathcal{H}\mathcal{R}$ (+CRF)	70.5	80.3	60.7

**Sequential decoder helps !**

## Results I

---

### Decoder Choice

	Avg	Avg DA	Avg E/S
BERT (+MLP)	72,8	81.5	64.0
BERT (+GRU)	69.9	80.4	59.3
BERT (+CRF)	72.8	81.5	64.1
$\mathcal{H}\mathcal{R}$ (+MLP)	69.8	79.1	60.4
$\mathcal{H}\mathcal{R}$ (+GRU)	67.6	79.4	55.7
$\mathcal{H}\mathcal{R}$ (+CRF)	70.5	80.3	60.7

### Convergence Speed

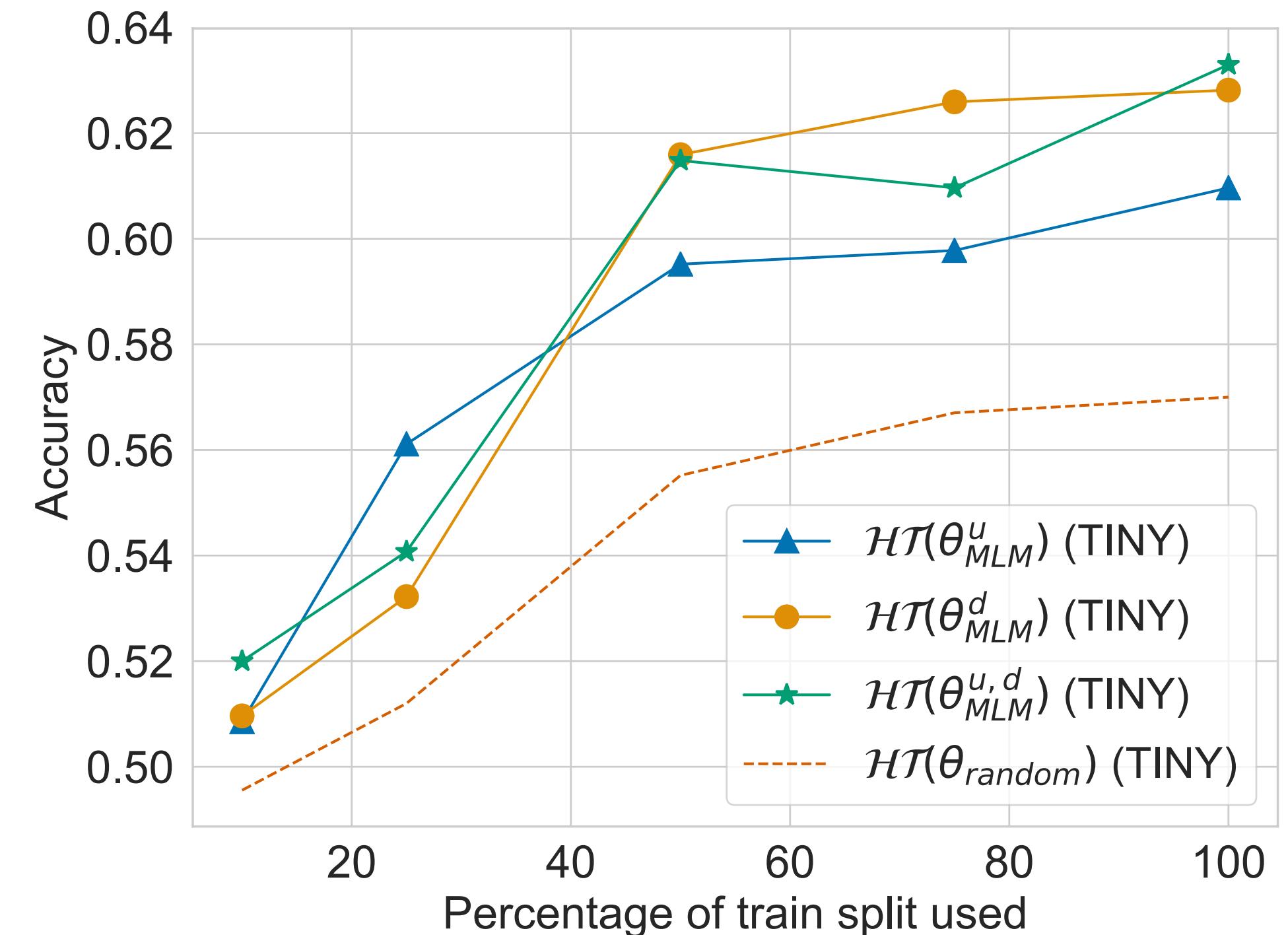
Sequential decoder helps !

## Results I

### Decoder Choice

	Avg	Avg DA	Avg E/S
BERT (+MLP)	72,8	81.5	64.0
BERT (+GRU)	69.9	80.4	59.3
BERT (+CRF)	72.8	81.5	64.1
$\mathcal{H}\mathcal{R}$ (+MLP)	69.8	79.1	60.4
$\mathcal{H}\mathcal{R}$ (+GRU)	67.6	79.4	55.7
$\mathcal{H}\mathcal{R}$ (+CRF)	70.5	80.3	60.7

### Convergence Speed



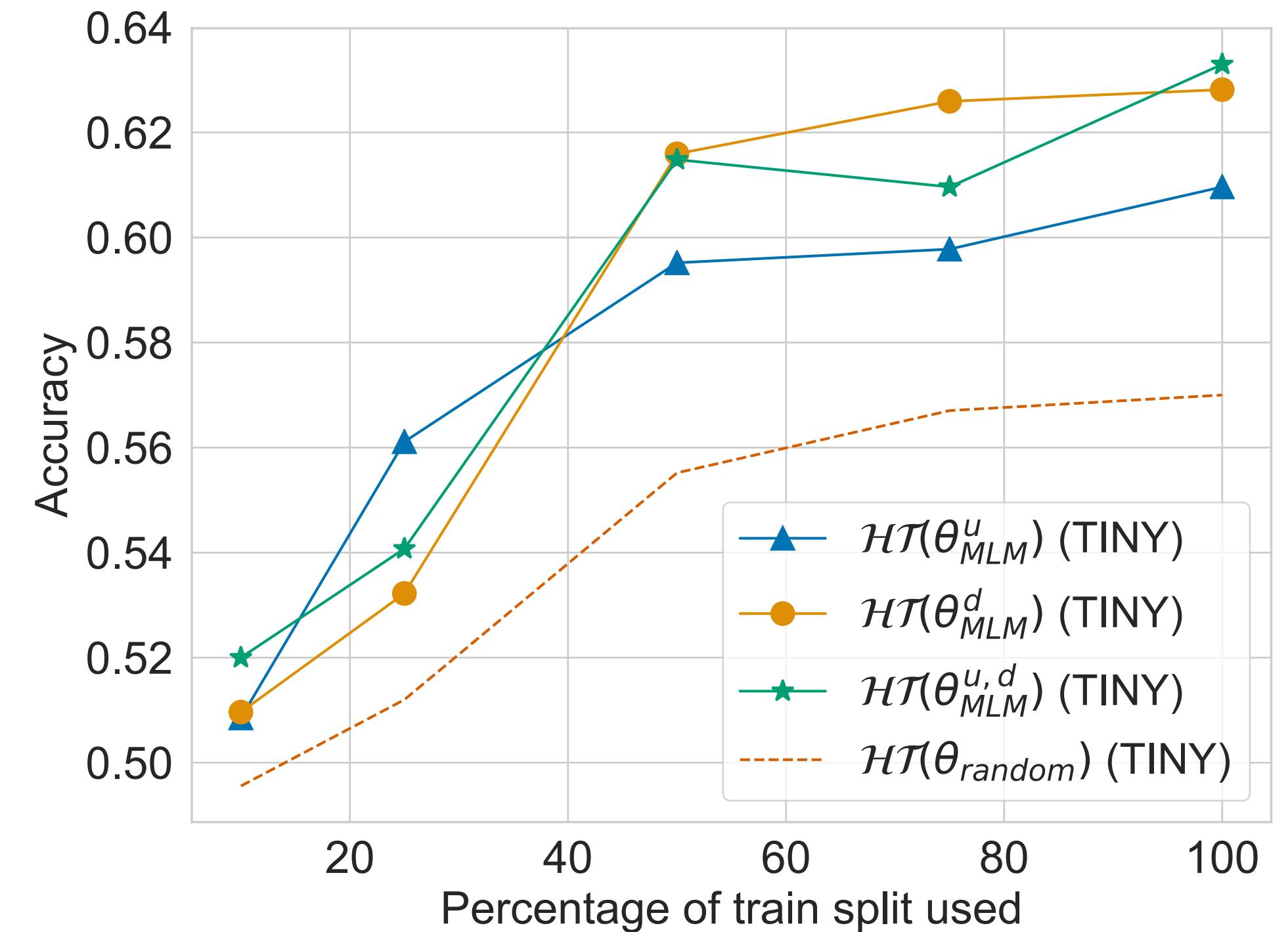
Sequential decoder helps !

## Results I

### Decoder Choice

	Avg	Avg DA	Avg E/S
BERT (+MLP)	72,8	81.5	64.0
BERT (+GRU)	69.9	80.4	59.3
BERT (+CRF)	72.8	81.5	64.1
$\mathcal{H}\mathcal{R}$ (+MLP)	69.8	79.1	60.4
$\mathcal{H}\mathcal{R}$ (+GRU)	67.6	79.4	55.7
$\mathcal{H}\mathcal{R}$ (+CRF)	70.5	80.3	60.7

### Convergence Speed



Sequential decoder helps !

Pretrained Representations  
Learn faster!

# Conclusion

---

## What is before?

**Guiding attention in Sequence-to-sequence models for Dialogue Act prediction**

Pierre Colombo\*, Emile Chapuis\*, Matteo Manica, Emmanuel Vignon, Giovanna Varni, and Chloé Clavel at AAAI 2020

## What is next?

**Code-switched inspired losses for generic spoken dialog representations**

Emile Chapuis\*, Pierre Colombo\*, Matthieu Labeau, Chloé Clavel at EMNLP 2021

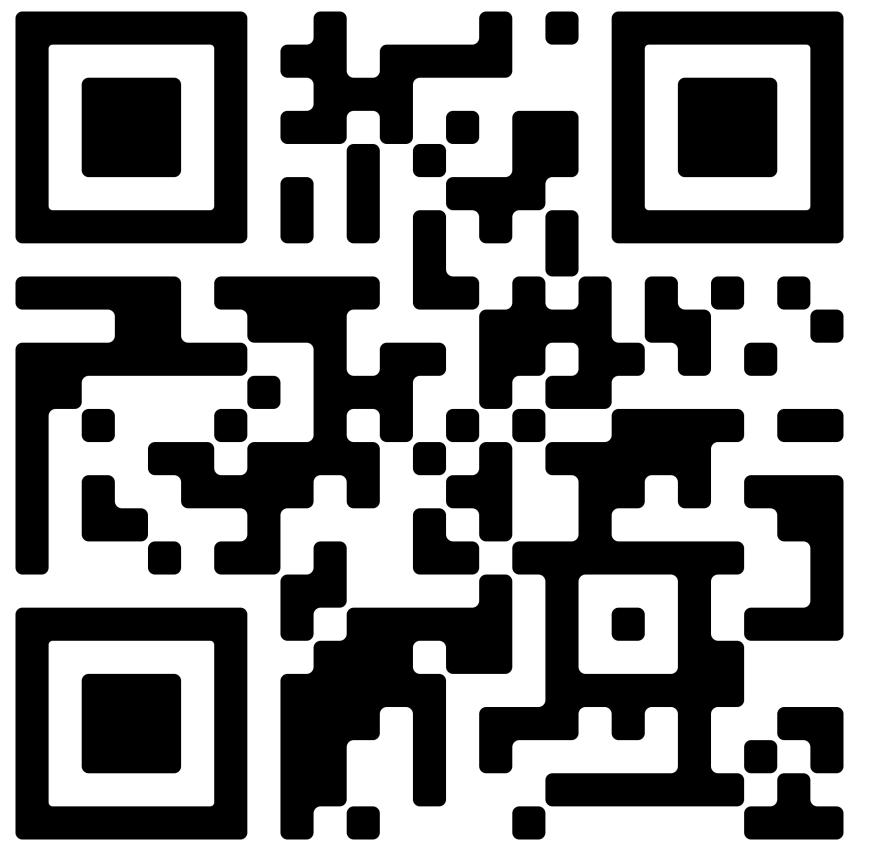
**Thesis of Emile Chapuis: Methods for Spoken Dialogue Understanding**

Décembre 2021

# Thanks for listening

Title: Hierarchical Pre-training for Sequence Labelling  
in SpokenDialog

[Link to Paper](#)



Corresponding Authors:

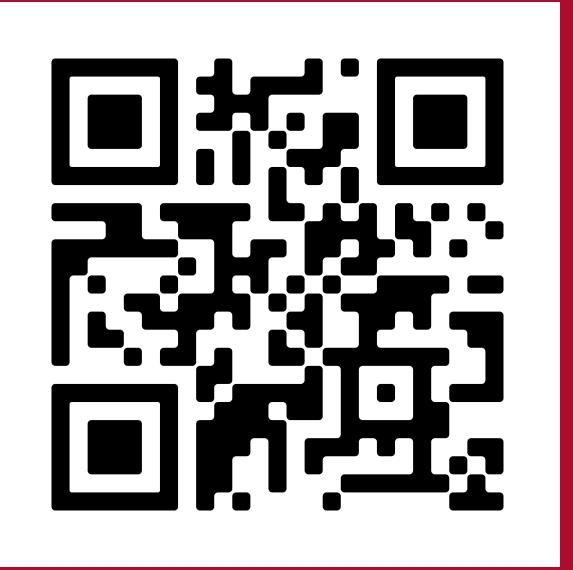


Emile Chapuis



Pierre Colombo

[Link to SILICONE](#)



# Context

---

# Context

---

## Conversational AI or dialogue systems

# Context

---

## Conversational AI or dialogue systems

- Natural Language Processing (**NLP**)
- Linguistics
- Psychology
- Information Retrieval (**IR**)
- Machine Learning (**ML**)

# Context

---

## Conversational AI or dialogue systems

- Natural Language Processing (**NLP**)
- Linguistics
- Psychology
- Information Retrieval (**IR**)
- Machine Learning (**ML**)



# Context

---

## Conversational AI or dialogue systems

- Natural Language Processing (**NLP**)
- Linguistics
- Psychology
- Information Retrieval (**IR**)
- Machine Learning (**ML**)



## Industrial applications

**Best available solutions are far from being capable of carrying open-domain conversations**

## Results II

---

## Results II

---

### Written vs Spoken

## Results II

---

### Written vs Spoken

	Avg DA	Avg E/S
BERT (4 layers)	80.5	60.2
$\mathcal{HT}(\theta_{BERT-2layers})$	80.5	61.1
$\mathcal{HT}(\theta_{MLM}^u)$	<b>80.8</b>	<b>64.0</b>

## Results II

---

### Written vs Spoken

	Avg DA	Avg E/S
BERT (4 layers)	80.5	60.2
$\mathcal{HT}(\theta_{BERT-2layers})$	80.5	61.1
$\mathcal{HT}(\theta_{MLM}^u)$	<b>80.8</b>	<b>64.0</b>

Training on Spoken Corpora help!

## Results II

---

### Written vs Spoken

	Avg DA	Avg E/S
BERT (4 layers)	80.5	60.2
$\mathcal{HT}(\theta_{BERT-2layers})$	80.5	61.1
$\mathcal{HT}(\theta_{MLM}^u)$	<b>80.8</b>	<b>64.0</b>

### Size Matters

Training on Spoken Corpora help!

## Results II

---

### Written vs Spoken

	Avg DA	Avg E/S
BERT (4 layers)	80.5	60.2
$\mathcal{HT}(\theta_{BERT-2layers})$	80.5	61.1
$\mathcal{HT}(\theta_{MLM}^u)$	<b>80.8</b>	<b>64.0</b>

### Size Matters

	Emb.	Word	Seq	Total
BERT	5*23	87		110
BERT (4-layer)		43		66
HMLP (TINY)		8.6	7.8	40
		2.9	2.8	28.7
(SMALL)		10.6	10.6	45

Training on Spoken Corpora help!

## Results II

---

### Written vs Spoken

	Avg DA	Avg E/S
BERT (4 layers)	80.5	60.2
$\mathcal{HT}(\theta_{BERT-2layers})$	80.5	61.1
$\mathcal{HT}(\theta_{MLM}^u)$	<b>80.8</b>	<b>64.0</b>

### Size Matters

	Emb.	Word	Seq	Total
BERT	5*23	87		110
BERT (4-layer)		43		66
HMLP (TINY)		8.6	7.8	40
		2.9	2.8	28.7
(SMALL)		10.6	10.6	45

Training on Spoken Corpora help!

Hierarchy factorize parameters:

- Less GPUs
- Faster Models
- Easier to learn

# Utterance Level GAP Loss

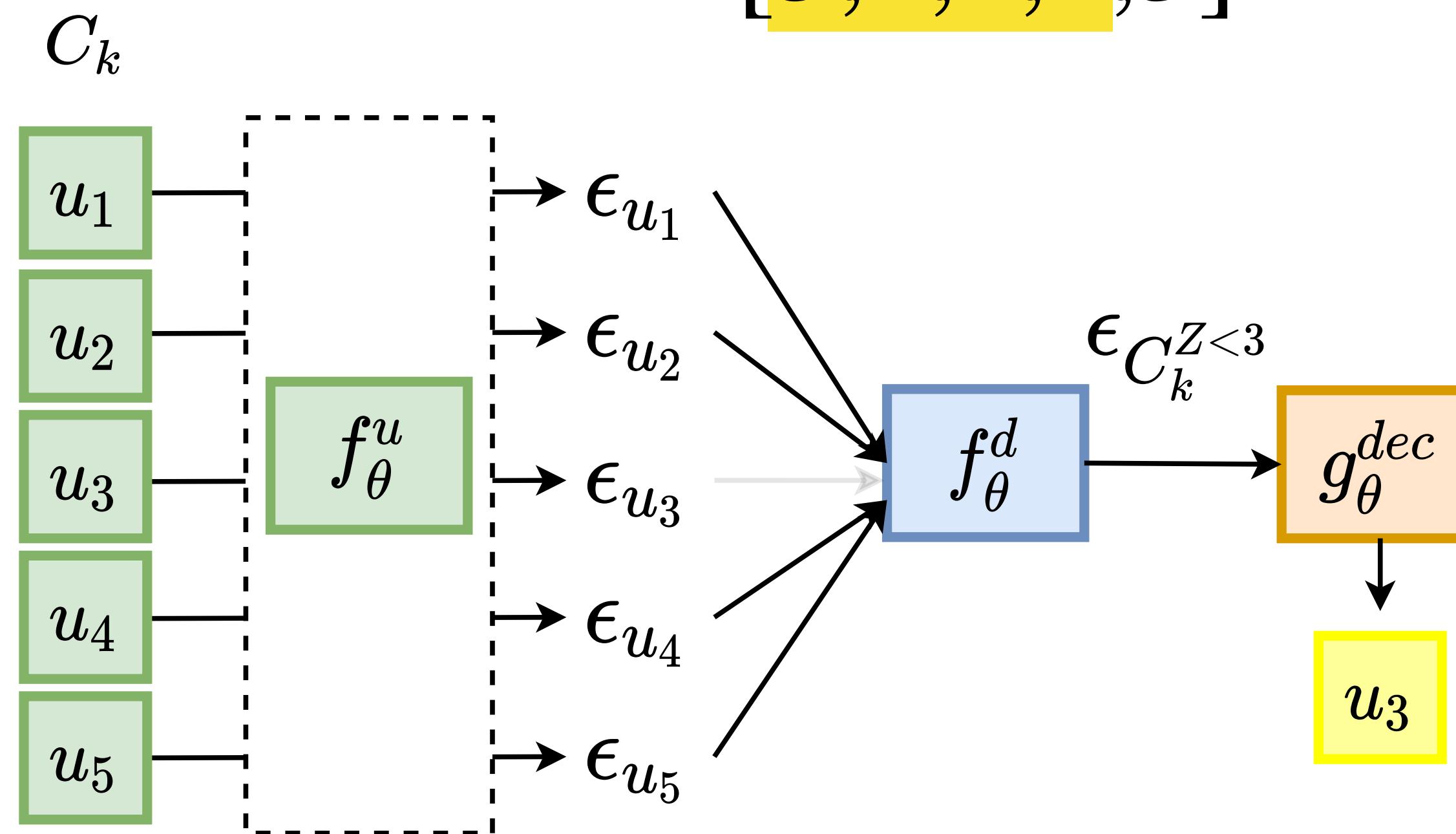
Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

Generalized Autoregressive Pre-training (GAP)

$$\mathcal{L}_{\text{GAP}}^d(\theta, C_k) = \mathbb{E}_{\mathbf{z} \sim \mathbb{Z}_T} \left[ \sum_{t=1}^{|C_k|} \sum_{i=1}^{|u_{z_t}|} \log p_\theta(\omega_i^{z_t} | C_k^{\mathbf{z} < t}) \right]$$

$$\mathbf{Z} = [5, 2, 4, 1, 3]$$



# Utterance Level GAP Loss

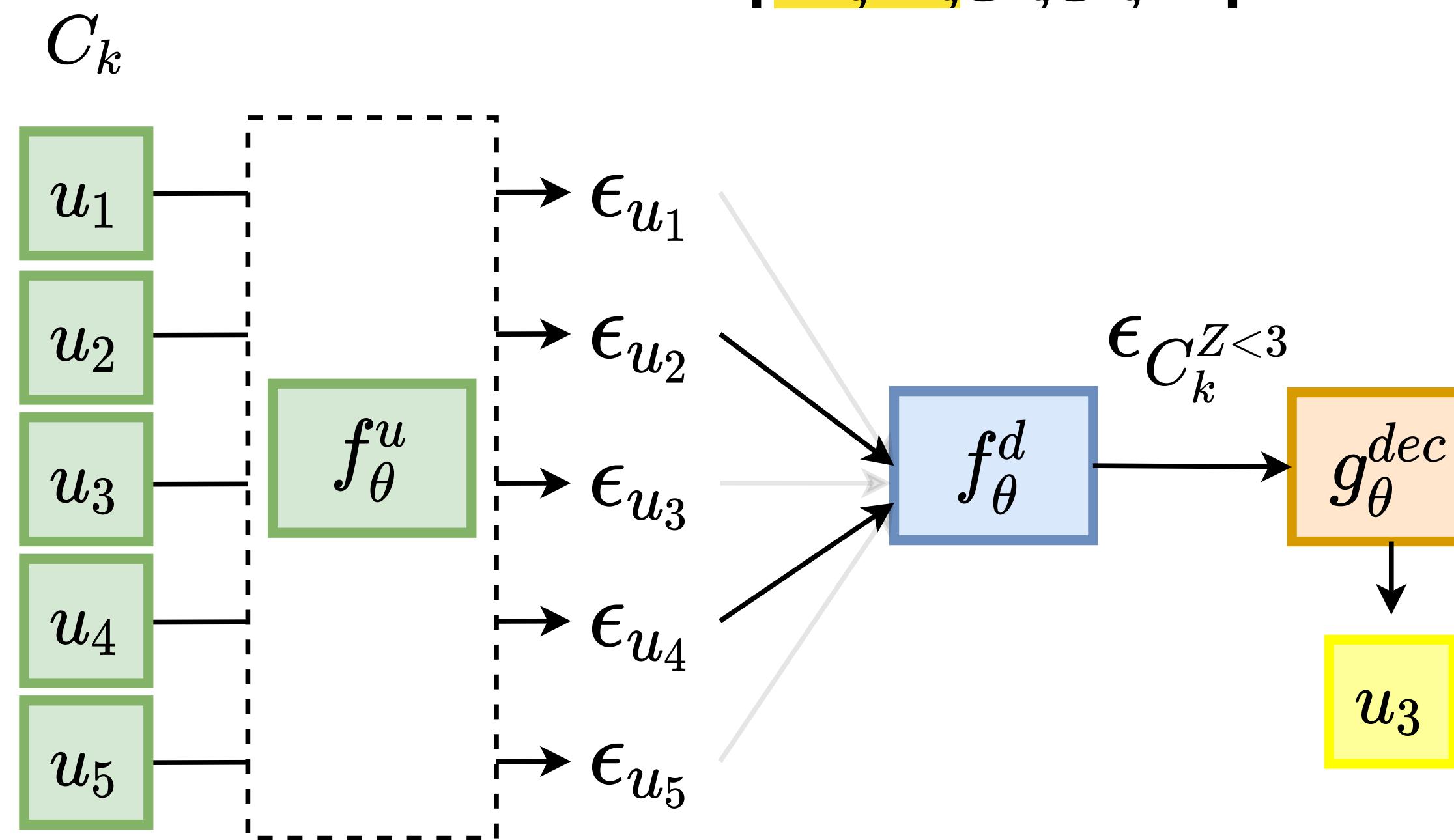
Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

Masked Language Model Loss  
(MLM)

$$\mathcal{L}_{\text{GAP}}^d(\theta, C_k) = \mathbb{E}_{\mathbf{z} \sim \mathbb{Z}_T} \left[ \sum_{t=1}^{|C_k|} \sum_{i=1}^{|u_{z_t}|} \log p_\theta(\omega_i^{z_t} | C_k^{\mathbf{z} < t}) \right]$$

$$\mathbf{Z} = [2, 4, 3, 5, 1]$$



# Utterance Level GAP Loss

Dialog Level Pretraining

$$\mathcal{L}(\theta) = \lambda_u \mathcal{L}^u(\theta) + \lambda_d \mathcal{L}^d(\theta).$$

Masked Language Model Loss  
(MLM)

$$\mathcal{L}_{\text{GAP}}^d(\theta, C_k) = \mathbb{E}_{\mathbf{z} \sim \mathbb{Z}_T} \left[ \sum_{t=1}^{|C_k|} \sum_{i=1}^{|u_{z_t}|} \log p_\theta(\omega_i^{z_t} | C_k^{\mathbf{z} < t}) \right]$$

$$\mathbf{Z} = [4, 3, 2, 5, 1]$$

