

## CCAI 2020 大会特邀报告(II)

微软亚洲研究院：

## 多语种和多模态任务的预训练模型

临菲信息技术港

CCAI 2020 第二天大会特邀报告。ppt 根据视频摘录。

## Pre-trained Models in Multi-Lingual and Multi-Modality Tasks

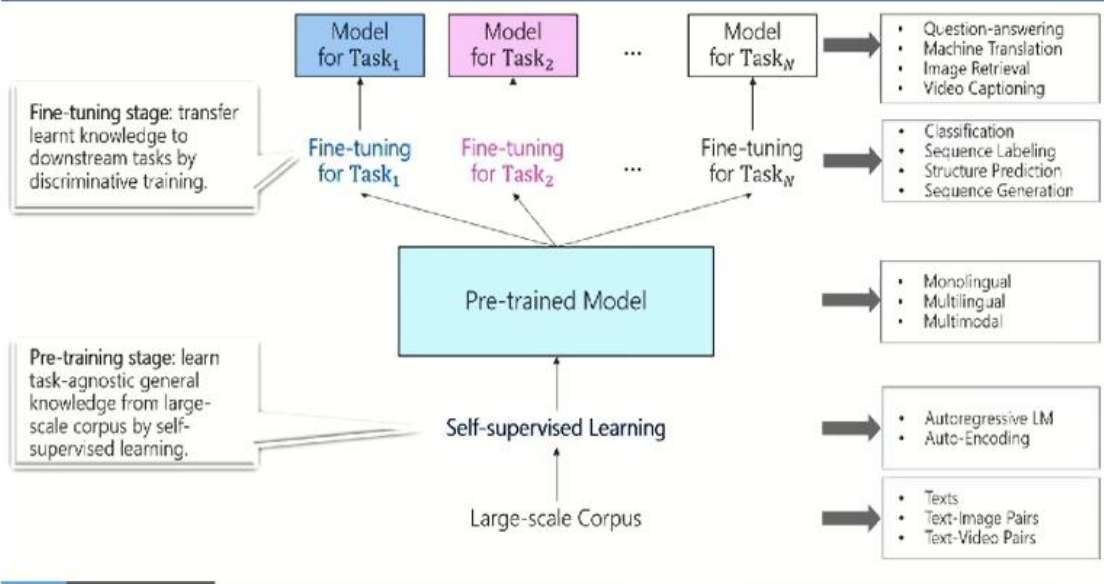
报告人：周明，微软亚洲研究院副院长。



# Agenda

1. Pre-trained models
2. Pre-trained models in multi-lingual tasks
3. Pre-trained models in multi-modality tasks
4. Summary

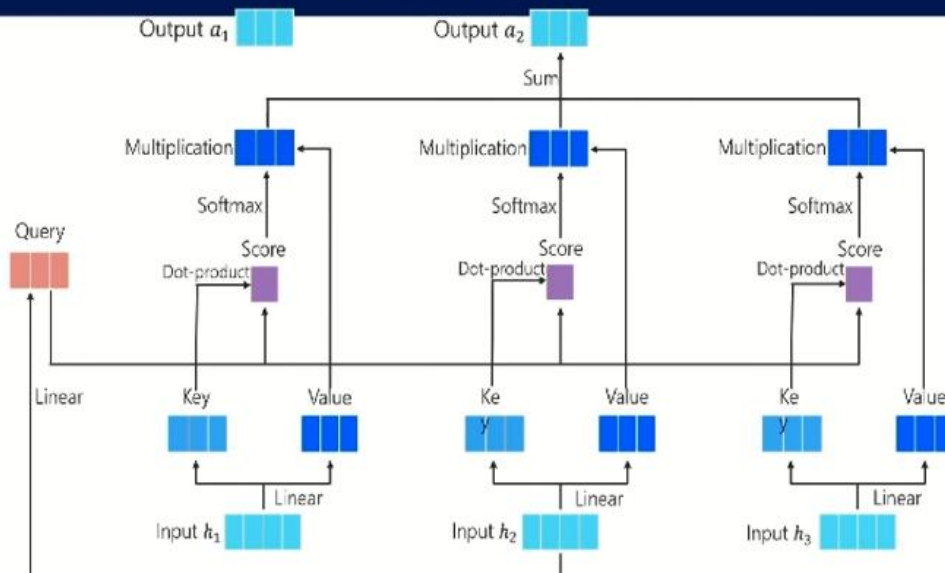
## Pre-trained Model: A New Paradigm of NLP



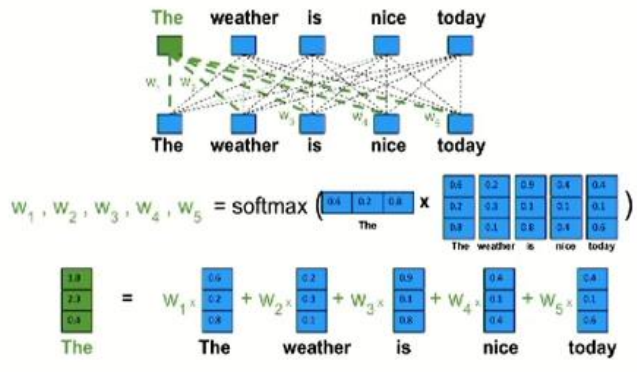
## Why Pre-trained Models?

1. Pre-trained models embed task-agnostic general knowledge.
  - Syntactic knowledge and semantic knowledge are implicitly encoded
2. Pre-trained models transfer learnt knowledge to downstream tasks.
  - Tasks with low-resource annotations tasks and languages
3. Pre-trained models supports almost all NLP tasks with SOTA results
  - Natural language understanding tasks and generation tasks
4. Provide a scalable solution to various applications
  - Supporting a new task only needs finetuning via task-specific labelled data

## Self-Attention

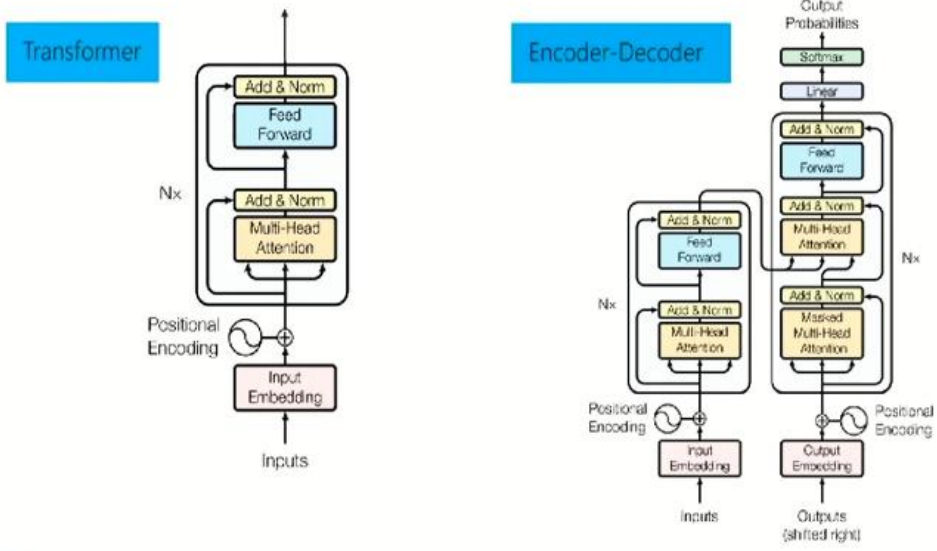


# A Simplified Example for Self-Attention



A simplified example of self-attention in Transformer

# Key Technologies (1): Transformer as Backbone



# Multi-Head Attention

**Multi-Head Attention**

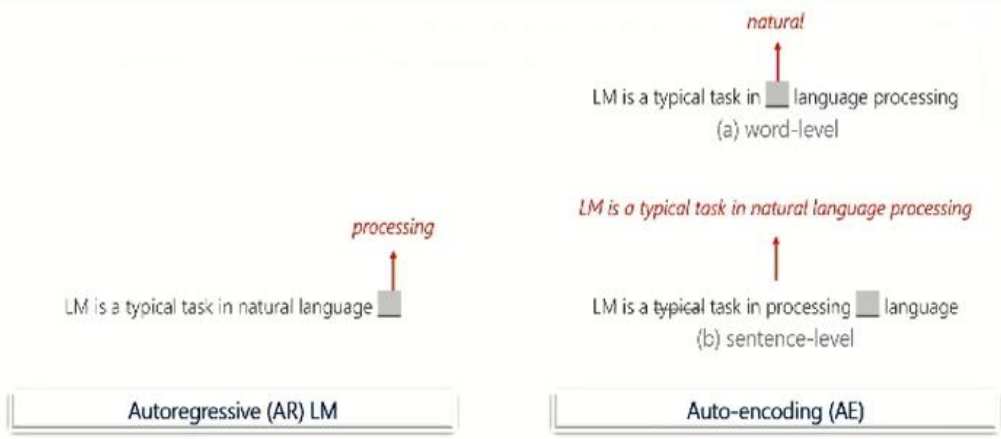
$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

- Multi-head attention allows the model to jointly attend to **information from different representation subspaces** at different positions.

<https://arxiv.org/pdf/1706.03762.pdf>

# Key Technologies (2): Pre-training by Self-supervised Learning



Self-supervised learning is a form of unsupervised learning where the data itself provides the supervision.

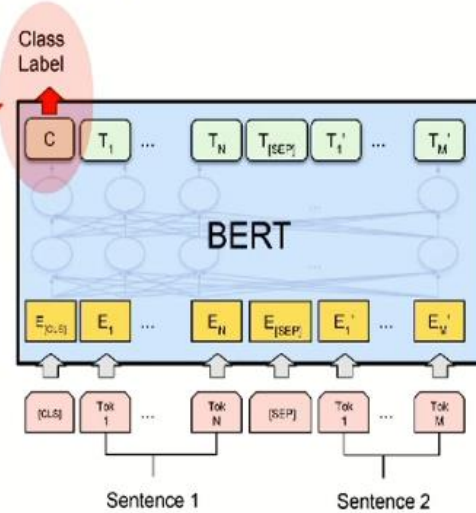
## Key Technologies (3): Fine-tuning by Discriminative Training

(take BERT-based Sentence Pair Matching as an example)

Given the final hidden vector  $C \in \mathbb{R}^H$  of the first input token ([CLS]), fine-tune BERT by a standard classification loss with  $C$  and  $W$ :

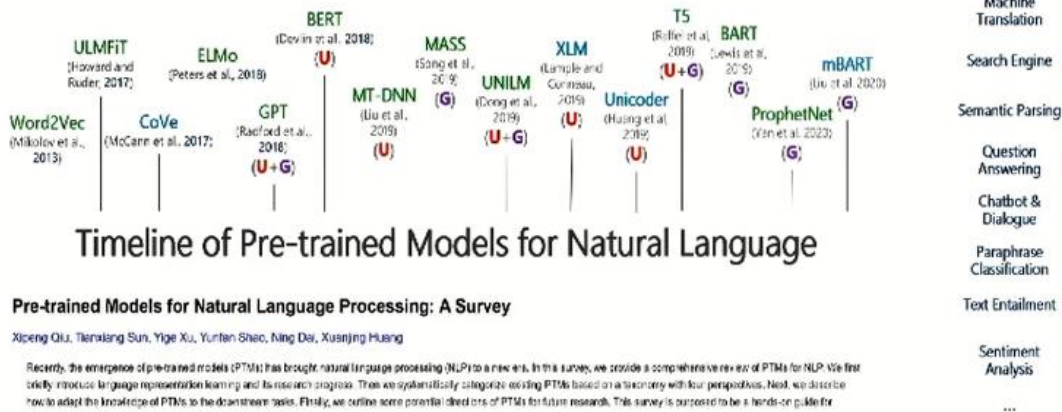
$$\log(\text{softmax}(CW^T))$$

where  $W \in \mathbb{R}^{K \times H}$  is a classification layer,  $K$  is the number of labels.



**GREEN:** monolingual pre-trained models  
**BLUE:** multilingual pre-trained models  
**U:** for understanding tasks  
**G:** for generation tasks

## Roadmap



### Pre-trained Models for Natural Language Processing: A Survey

Xicong Qiu, Tianyang Sun, Yige Xu, Yunfeng Shao, Ning Dai, Xuanjing Huang

Recently, the emergence of pre-trained models (PTMs) has brought natural language processing (NLP) to a new era. In this survey, we provide a comprehensive review of PTMs for NLP. We first briefly introduce language representation learning and its research progress. Then we systematically categorize existing PTMs based on a taxonomy with four perspectives. Next, we describe how to adapt the knowledge of PTMs to the downstream tasks. Finally, we outline some potential directions of PTMs for future research. This survey is supposed to be a hands-on guide for understanding, using, and developing PTMs for various NLP tasks.

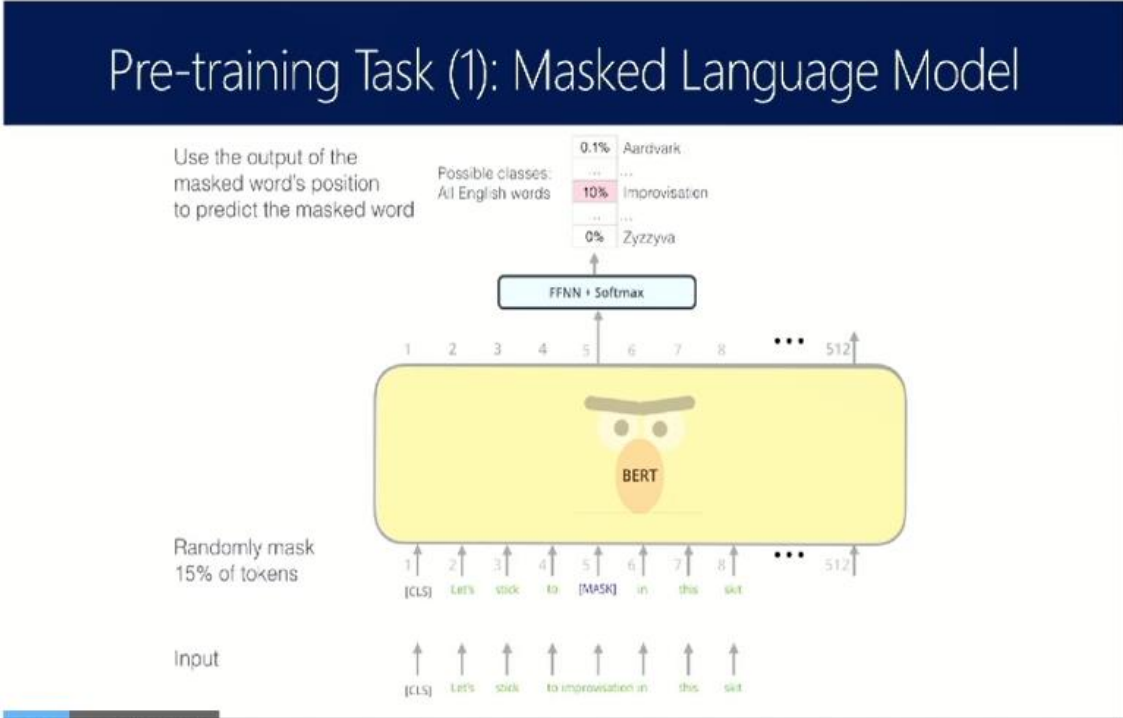
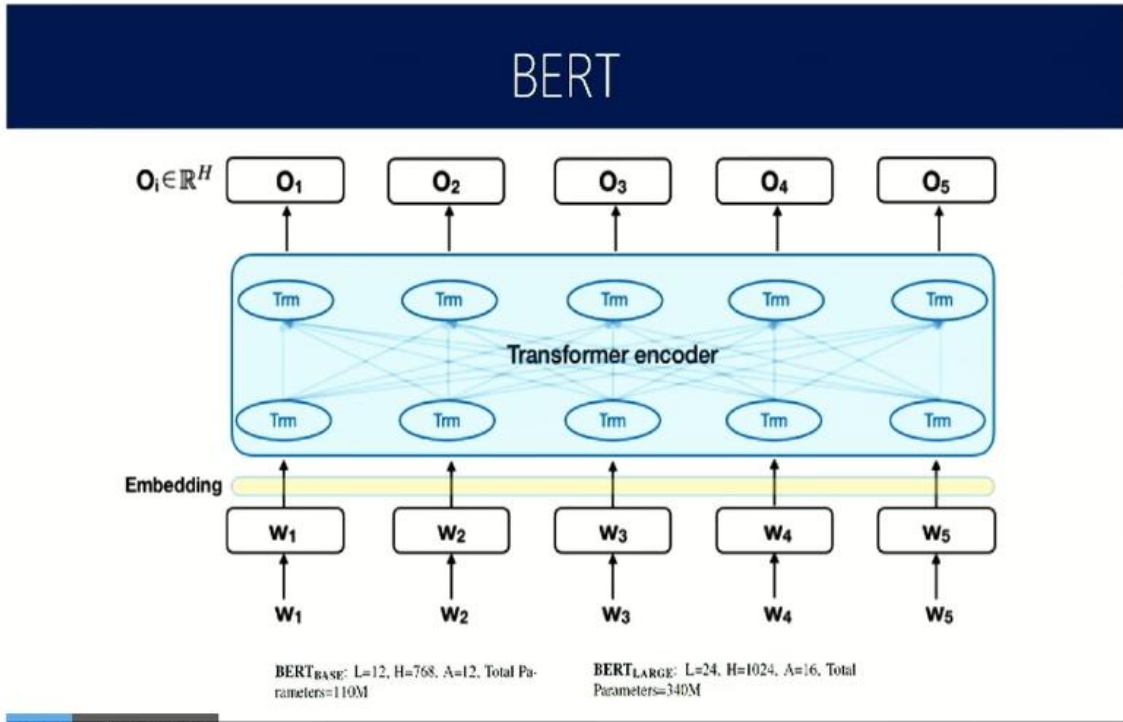
<https://arxiv.org/abs/2003.08271>

## Connections and Differences Between (Monolingual) Pre-trained Models

Model Name	Model Usage	Model Backbone	Model Contribution
GPT (OpenAI)	Understanding & Generation	Transformer Encoder	1 <sup>st</sup> unidirectional pre-trained LM based on Transformer
BERT (Google)	Understanding	Transformer Encoder	1 <sup>st</sup> bidirectional pre-trained LM based on Transformer
MT-DNN (MS)	Understanding	Transformer Encoder	use multiple understanding tasks in pre-training
MASS (MS)	Generation	Separate Transformer Encoder-Decoder	use masked span prediction for generation tasks
UniLM (MS)	Understanding & Generation	Unified Transformer Encoder-Decoder	unify understanding and generation tasks in pre-training with different attention masks
RoBERTa (FB)	Understanding	Transformer Encoder	use better pre-training tricks, such as dynamic masking, large batches, removing NSP, data sampling
ERNIE (Baidu)	Understanding	Transformer Encoder	prove noun phrase masking and entity masking are better than word masking
SpanBERT (FB)	Understanding	Transformer Encoder	prove random span masking is better than others
XLNet (Google)	Understanding	Transformer Encoder	unify autoregressive LM and autoencoding tasks in pre-training with the two-stream self-attention
T5 (Google)	Generation	Separate Transformer Encoder-Decoder	use a separate encoder-decoder for understanding and generation tasks and prove it is the best choice; compare different hyper-parameters and show the best settings
BART (FB)	Generation	Separate Transformer Encoder-Decoder	try different text noising methods for generation tasks
ELECTRA (Google)	Understanding	Transformer Generator-Discriminator	use a simple but effective GAN-style pre-training task
ProphetNet (MS)	Generation	Separate Transformer Encoder-Decoder	use future n-gram prediction for generation tasks with the n-stream self-attention

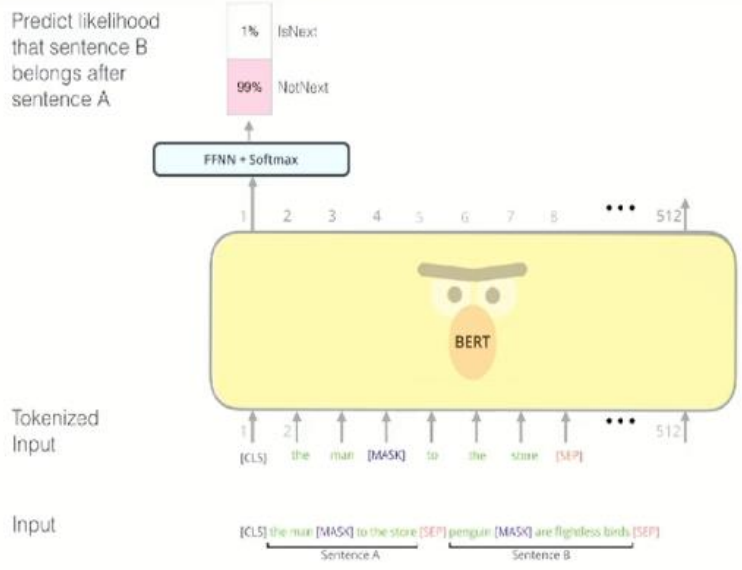
## Push the Boundary of Pre-trained Models

- Model size (extremely big models, e.g. Turing 17B, GPT-3 175B parameters...)
- Pre-training methods and models (pre-training tasks, masking policy, network structures ...)
- From single language, to multi-lingual, to multi-modal
- Model compression, knowledge distilling, economical models for practical needs





# Pre-training Task (2): Next Sentence Prediction



# UniLM (Dong et al., 2019)

Allow to attend  
 Prevent from attending

**Unified LM with Shared Parameters**

**Self-attention Masks**

	ROUGE-1	ROUGE-2	ROUGE-L
<i>Extractive Summarization</i>			
LEAD-3	40.42	17.62	36.67
Best Extractive	43.25	20.24	39.63
<i>Abstractive Summarization</i>			
PGNet	39.53	17.28	37.98
Bottom-Up	41.22	18.68	38.34
S2S-ELMo	41.56	18.94	38.47
UNILM	<b>43.47</b>	<b>20.30</b>	<b>40.63</b>

Table 3: Evaluation results on CNN/DailyMail. Models in the first block are extractive systems listed here for reference, while the others are abstractive models. The results of the best reported extractive model are taken from (Liu, 2019).

	EM	F1
RMR+ELMo (Hu et al., 2018)	71.4	73.7
BERT <sub>Large</sub>	78.9	81.8
UNILM	<b>80.5</b>	<b>83.4</b>

Table 4: Extractive question answering results on the SQuAD development set.

# Application: Question Answering

when were women allowed to vote in the usa

1920 (100 RESULTS) Article

Ratified on **August 18, 1920**, the 19th Amendment to the U.S. Constitution granted American women the right to vote—a right known as woman suffrage.

19th Amendment - Women's History - HISTORY.com

Women's suffrage in the United States - Wikipedia

19TH AMENDMENT

Ratified on August 18, 1920, the 19th Amendment to the U.S. Constitution granted American women the right to vote—a right known as woman suffrage.

# Application: Question Generation

## Text Passage

in **1066**<sup>1,2</sup>, **duke william ii**<sup>3</sup> of normandy conquered england killing king harold ii at the battle of hastings. **the invading normans and their descendants**<sup>4</sup> replaced the anglo-saxons as the ruling class of england.

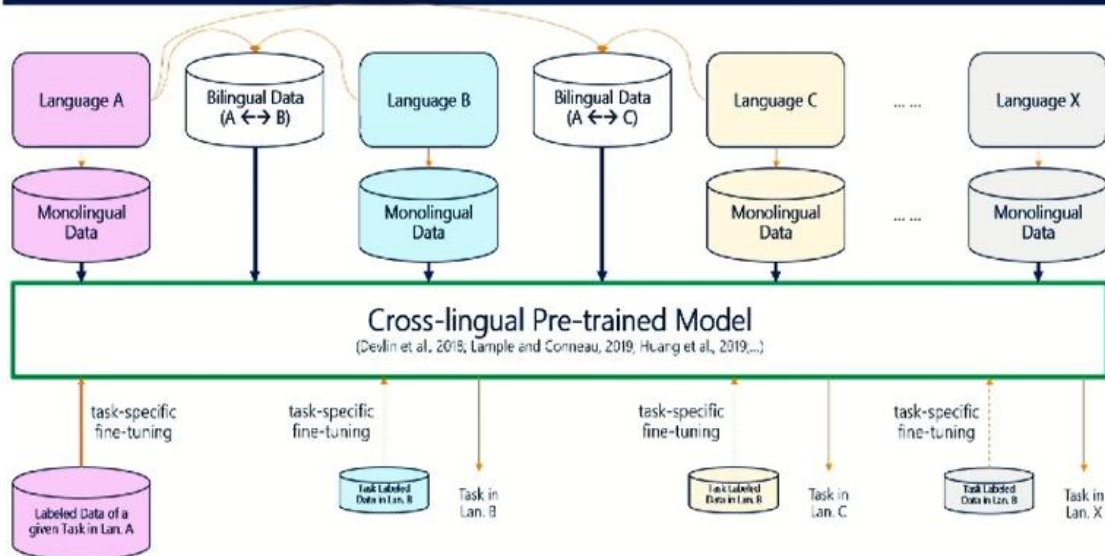
## Questions Generated by our System

- 1) when did the battle of hastings take place?
- 2) in what year was the battle of hastings fought?
- 3) who conquered king harold ii at the battle of hastings?
- 4) who became the ruling class of england?

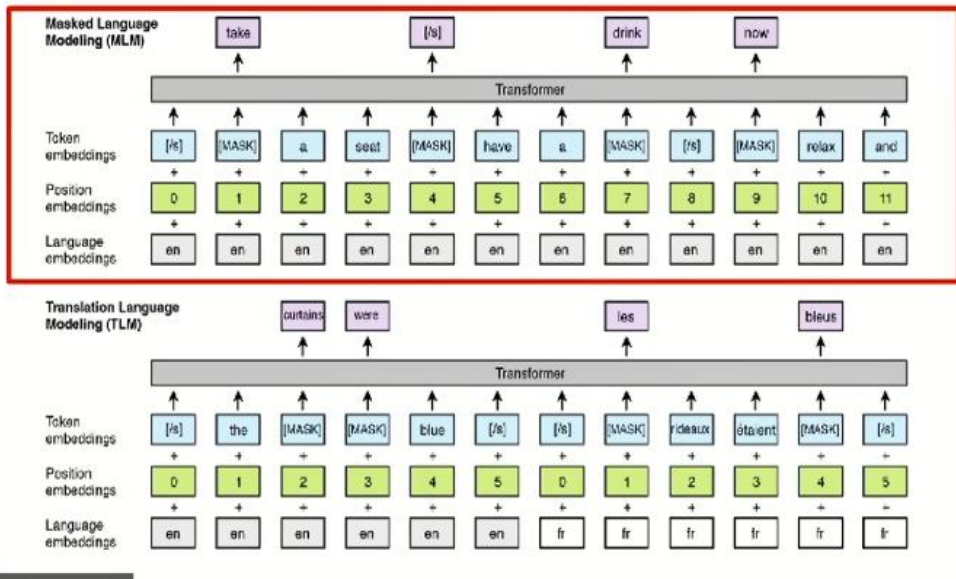
## Agenda

1. Pre-trained models
2. Pre-trained models in multi-lingual tasks
3. Pre-trained models in multi-modality tasks
4. Summary

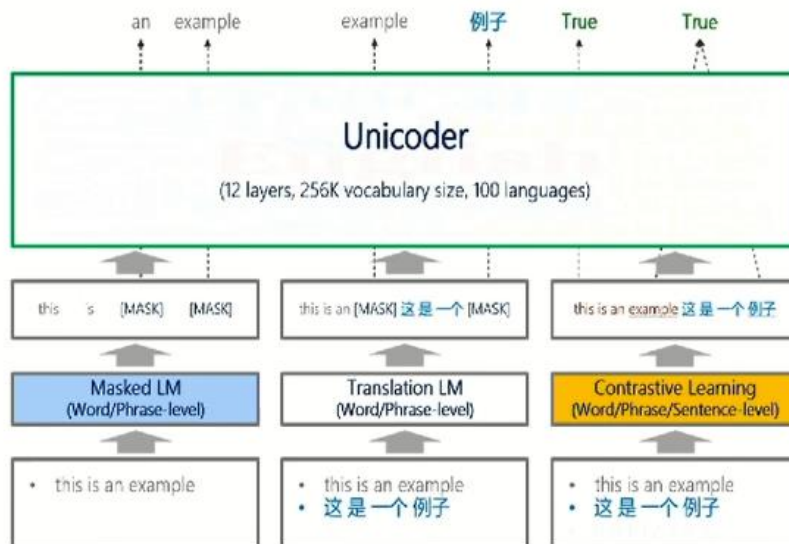
## Cross-Language Pre-training



## XLM (Lample and Conneau, 2019; Conneau et al., 2019)

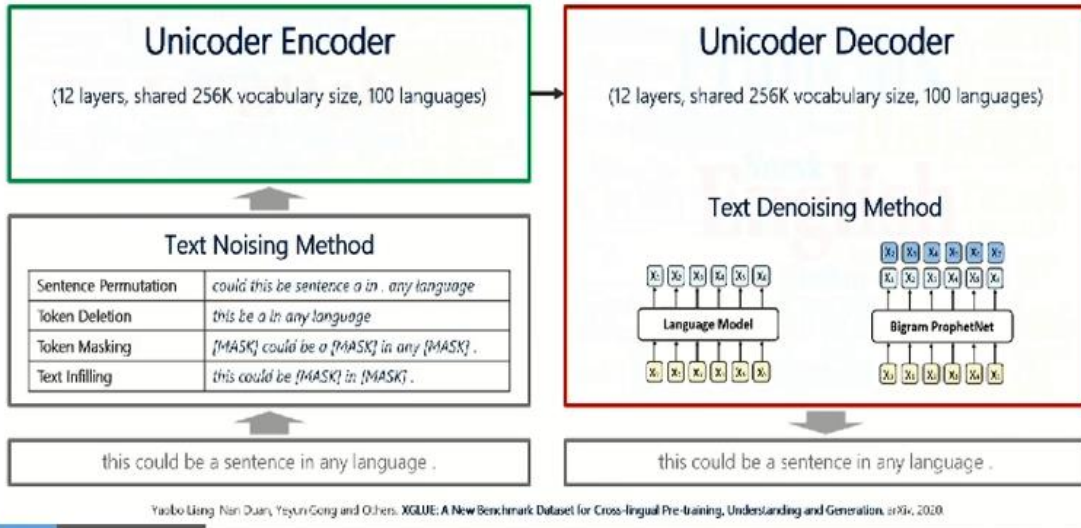


## Unicoder-1 (Huang et al., 2019)



Huawei Huang, Yanbo Liang, Nian Duan, Ming Gong, Lijun Shou, Dazhi Fang, Ming Zhou. Unicoder: A Universal Language Encoder by Pre-training with Multiple Cross-lingual Tasks. EMNLP, 2019

# Unicoder-2 (Liang et al., 2020)



## Tasks in XGLUE

Task	# of Languages	[Train] <sup>en</sup>	[Dev] <sup>en</sup>	[Test] <sup>en</sup>	Metric	Data Source
NER	4	150K	2.8K	3.4K	F1	ECI Multilingual Text Corpus
POS	18	25.4K	1.0K	0.9K	ACC	UD Tree-banks (v2.5)
NC*	5	100K	10K	10K	ACC	Commercial News Website
MLQA	7	87.6K	0.6K	5.7K	F1	Wikipedia
XNLI	15	433K	2.5K	5K	ACC	MultiNLI Corpus
PAWS-X	4	49.4K	2K	2K	ACC	Wikipedia
QADSM*	3	100K	10K	10K	ACC	Commercial Search Engine
WPR*	7	100K	10K	10K	nDCG	Commercial Search Engine
QAM*	3	100K	10K	10K	ACC	Commercial Search Engine
QG*	6	100K	10K	10K	BLEU-4	Commercial Search Engine
NTG*	5	300K	10K	10K	BLEU-4	Commercial News Website

Table 2: 11 downstream tasks in XGLUE. For each task, training set is only available in English. [Train]<sup>en</sup> denotes the number of labeled instances in the training set. [Dev]<sup>en</sup> and [Test]<sup>en</sup> denote the average numbers of labeled instances in the dev sets and test sets, respectively. \* denotes the corresponding dataset is constructed by this paper.

Task	ar	bg	de	el	en	es	fr	hi	it	nl	pl	pt	ru	sw	th	tr	ur	vi	zh
NER			✓		✓	✓				✓									
POS	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
NC*			✓		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
MLQA	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
XNLI	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
PAWS-X			✓		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
QADSM*			✓		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
WPR*			✓		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
QAM*			✓		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
QG*			✓		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
NTG*			✓		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 3: The 19 languages covered by the 11 downstream tasks: Arabic (ar), Bulgarian (bg), German (de), Greek (el), English (en), Spanish (es), French (fr), Hindi (hi), Italian (it), Dutch (nl), Polish (pl), Portuguese (pt), Russian (ru), Swahili (sw), Thai (th), Turkish (tr), Urdu (ur), Vietnamese (vi), and Chinese (zh). All these 6 new tasks with \* are labeled by human, except es, it and pt datasets in QG (80+% accuracy) are obtained by an in-house QA ranker.

- NER: Named Entity Recognition
- POS: Part-of-Speech Tagging
- NC: News Classification
- MLQA: Multilingual MRC
- XNLI: Natural Language Inference
- PAWS-X: Paraphrase Classification
- QADSM: Query-Ads Matching
- WPR: Web Page Ranking
- QAM: Question-Answer Matching
- QG: Question Generation
- NTG: News Title Generation

Task	Model	ar	bg	de	el	en	es	fr	hi	it	id	pt	pt	ru	sv	th	tr	ur	vi	zh	AVG
NER	M-BERT	-	-	69.2	-	90.6	75.4	-	-	73.9	-	-	-	-	-	-	-	-	-	-	78.2
	XLM-R <sub>base</sub>	-	-	70.4	-	90.9	75.2	-	-	70.5	-	-	-	-	-	-	-	-	-	-	79.0
	Uniswift <sub>0.1</sub>	-	-	71.8	-	91.1	74.1	-	-	81.6	-	-	-	-	-	-	-	-	-	-	76.7
POS	M-BERT	52.4	85.0	38.7	81.5	95.6	39.8	37.6	58.4	91.3	88.0	11.8	88.3	78.8	-	42.3	69.2	53.8	54.1	58.3	74.7
	XLM-R <sub>base</sub>	67.3	88.8	92.2	84.2	96.2	39.0	39.9	74.5	92.6	88.5	15.4	89.7	86.0	-	57.9	72.7	62.1	55.2	60.4	70.8
	Uniswift <sub>0.1</sub>	68.6	88.5	92.0	83.3	96.1	39.1	39.4	69.0	92.5	88.9	13.6	89.8	86.7	-	57.6	73.0	53.8	56.3	60.2	70.6
NC	M-BERT	-	-	82.0	-	92.2	81.6	79.0	-	-	-	-	-	79.0	-	-	-	-	-	-	82.7
	XLM-R <sub>base</sub>	-	-	84.5	-	91.8	81.2	78.2	-	-	-	-	-	79.4	-	-	-	-	-	-	83.4
	Uniswift <sub>0.1</sub>	-	-	84.2	-	91.7	83.5	78.5	-	-	-	-	-	79.7	-	-	-	-	-	-	83.5
MLQA	M-BERT	50.9	-	65.5	-	80.5	67.1	-	47.0	-	-	-	-	-	-	-	-	-	-	-	59.3
	XLM-R <sub>base</sub>	56.4	-	62.1	-	80.1	67.8	-	60.3	-	-	-	-	-	-	-	-	-	-	-	67.1
	Uniswift <sub>0.1</sub>	57.8	-	62.7	-	80.9	68.6	-	62.7	-	-	-	-	-	-	-	-	-	-	-	67.5
XNLI	M-BERT	64.0	68.9	71.1	66.4	82.1	74.3	71.8	60.0	-	-	-	-	69.0	50.4	55.8	61.6	53.0	60.5	65.3	65.3
	XLM	73.1	77.4	77.5	76.6	83.0	78.3	78.7	69.8	-	-	-	-	73.2	68.4	72.2	72.5	67.3	76.1	78.5	75.1
	Uniswift <sub>0.1</sub>	72.1	77.5	77.0	75.9	80.6	79.2	78.2	69.8	-	-	-	-	73.5	64.7	71.6	72.9	65.1	74.4	77.7	74.2
PAWS-X	M-BERT	-	-	82.0	-	94.0	85.9	85.0	-	-	-	-	-	-	-	-	-	-	-	-	87.2
	XLM-R <sub>base</sub>	-	-	89.9	-	94.4	88.0	89.7	-	-	-	-	-	-	-	-	-	-	-	-	90.3
	Uniswift <sub>0.1</sub>	-	-	87.0	-	94.9	88.8	89.3	-	-	-	-	-	-	-	-	-	-	-	-	90.1
QADSM	M-BERT	-	-	69.3	-	68.1	-	64.1	-	-	-	-	-	-	-	-	-	-	-	-	64.2
	XLM-R <sub>base</sub>	-	-	65.4	-	71.7	-	63.3	-	-	-	-	-	-	-	-	-	-	-	-	65.6
	Uniswift <sub>0.1</sub>	-	-	64.6	-	71.8	-	63.7	-	-	-	-	-	-	-	-	-	-	-	-	65.4
WTR	M-BERT	-	-	80.6	-	78.1	75.1	74.2	-	70.1	-	-	-	76.6	-	-	-	-	-	-	84.2
	XLM-R <sub>base</sub>	-	-	77.6	-	78.2	76.0	74.4	-	70.7	-	-	-	77.5	-	-	-	-	-	-	83.9
	Uniswift <sub>0.1</sub>	-	-	77.2	-	78.4	75.7	74.9	-	70.5	-	-	-	77.4	-	-	-	-	-	-	84.4
QAM	M-BERT	-	-	66.7	-	67.5	-	65.0	-	-	-	-	-	-	-	-	-	-	-	-	65.1
	XLM-R <sub>base</sub>	-	-	68.1	-	66.3	-	67.8	-	-	-	-	-	-	-	-	-	-	-	-	65.4
	Uniswift <sub>0.1</sub>	-	-	68.4	-	66.9	-	67.4	-	-	-	-	-	-	-	-	-	-	-	-	65.9
AVG <sub>U</sub>	M-BERT	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	72.6
	XLM-R <sub>base</sub>	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	75.4
	Uniswift <sub>0.1</sub>	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	76.3
QG	M-BERT	-	-	5.1	-	7.8	0.1	0.1	-	0.2	-	-	-	0.1	-	-	-	-	-	-	1.1
	XLM-R <sub>base</sub>	-	-	6.1	-	5.0	0.0	0.0	-	0.1	-	-	-	0.0	-	-	-	-	-	-	1.0
	Uniswift <sub>0.1</sub> FR	-	-	3.0	-	14.0	12.4	4.2	-	15.8	-	-	-	8.3	-	-	-	-	-	-	9.6
NTG	M-BERT	-	-	0.7	-	3.0	0.4	0.4	-	-	-	-	-	0.0	-	-	-	-	-	-	2.1
	XLM-R <sub>base</sub>	-	-	0.6	-	3.1	0.4	0.3	-	-	-	-	-	0.0	-	-	-	-	-	-	1.9
	Uniswift <sub>0.1</sub> FR	-	-	1.8	-	15.6	0.0	3.7	-	-	-	-	-	7.7	-	-	-	-	-	-	9.6
AVG <sub>U</sub>	M-BERT	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1.7
	XLM-R <sub>base</sub>	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1.5
	Uniswift <sub>0.1</sub> FR	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	9.6
																					10.7

## Application: Multilingual Question Answering

ou est la plus grand usine du sucre au monde

27 130 000 Results

La raffinerie de sucre du groupe algérien Cevital produit 2,7 millions de tonnes de sucre par an, ce qui en fait la raffinerie la plus importante au monde. Cette raffinerie a doublé ses exportations de sucre blanc en passant de 377 000 tonnes à 600 000 tonnes en 2012.

Industrie sucrière — Wikipédia

**Bing fr-FR**

*English translation: where is the largest sugar factory in the world*

*English translation: The sugar refinery of the Algerian group Cevital produces 2.7 million tonnes of sugar a year, making it the largest refinery in the world. This refinery doubled its exports of white sugar from 377,000 tonnes to 600,000 tonnes in 2012.*

browserververlauf löschen windows 10

5 220 000 Results

Im Reiter "Allgemein" finden Sie den Unterpunkt "Browserververlauf". Klicken Sie dort auf die Schaltfläche "Löschen...". Es öffnet sich ein Fenster mit dem Namen "Browserverlauf löschen". Setzen Sie die Häkchen bei den Daten, die gelöscht werden sollen. Klicken Sie auf "Löschen". Der Verlauf wurde nun entfernt.

Verlauf löschen: Anleitung für alle Browser - CHIL

**Bing de-DE**

*English translation: delete browser history windows 10*

*English translation: In the tab "General" you will find the sub-item "Browser History". Click on the "Delete ..." button there. A window with the name "delete browser history" will open. Check the data you want to delete. Click on "Delete". The history has now been removed.*

## Application: Multilingual News Headline Generation

en	Input News	if you're planning a trip to europe , you probably want to check some famous landmarks off your list . but there are certain tourist traps you're better off missing . susana victoria perez has more .
	Golden Title	do yourself a favor and avoid these tourist traps in europe
	Unicoder <sup>DAB</sup> <sub>ZC</sub>	tourist traps you should avoid in europe
fr	Input News	alain juppe , candidat a la primaire de la droite , " ne se sent pas engage " par les investitures decidees par le parti les republicains preside par nicolas sarkozy , a affirme jeudi a l'ailp son directeur de campagne , gilles boyer . " c' est un processus mene a la hussarde , il n' y a pas de volonte d' equilibre et de rassemblement " , a-t-il denonce , en affirmant que " l' accord politique " entre les differents candidats a la primaire " n' a pas ete respecte " .
	Golden Title	legislatives : juppe " ne se sent pas engage " par les investitures
	Unicoder <sup>DAB</sup> <sub>ZC</sub>	alain juppe : " ne se sent pas engage " par les investitures
de	Input News	vermutlich zur verteidigung seines reviers hat ein aggressiver bussard in baden-wuerttemberg einen radfahrer zu fall gebracht . der sich dabei schwer verletzte . wie die polizei in ludwigsbug am freitag mitteilte , attackierte der greifvogel den 51-jahrigen am vortag auf einem radweg entlang einer landesstrasse . der bussard flog demnach so tief auf den radler zu , dass dieser ausweichen musste und sturzte . den angaben zufolge erlitt der mann schwere verletzungen und wurde von rettungskraefen in ein krankenhaus gebracht . " aus luftiger hoehe , von einem laternenmast aus , beobachtete der raubvogel anschliessend die unfallaufnahme " , hielt es im polizeibericht .
	Golden Title	aggressiver bussard bringt radfahrer zu fall
	Unicoder <sup>DAB</sup> <sub>ZC</sub>	aggressiver bussard in ludwigsbug stuerzes radler
es	Input News	despues de la marcha de bruce willis por problemas de agenda , steve carrell le sustituirá asi en la nueva pelicula que prepara woody allen . segun informa variety , el actor se une al reparto ya formado por blake lively , parker posey , kristen stewart , jesse eisenberg , jeannie berlin , corey stoll , anna camp , y ken stott , entre otros . como siempre , los detalles de la trama son aun un secreto aunque el rodaje se encuentre actualmente en marcha . por otro lado , aun no hay fecha de estreno ni distribuidora para la pelicula sin titulo de woody allen . sin embargo , el director tiene aun pendiente de estreno su ultimo filme con emma stone y joaquin phoenix titula da irrational man que se estrenara el proximo 25 de septiembre .
	Golden Title	steve carrell sustituye a bruce willis en la nueva pelicula de woody allen
	Unicoder <sup>DAB</sup> <sub>ZC</sub>	steve carrell sustituirá a steve carrell en woody allen

## Summary of Multilingual Pre-trained Models

Multilingual pre-trained models alleviate the low-resource issue in multiple languages

- Fine-tune on rich-resource languages (such as EN) and then apply to low-resource languages
- Show strong cross-lingual transfer capabilities in zero-shot or few-shot settings
- Achieve state-of-the-art results on various cross-lingual benchmarks.

Multilingual pretrained models can benefit many real-world applications

- Multilingual Search/QA/Ads/News/Text Summarization/...
- Low-resource Neural Machine Translation

Multilingual pretrained models still have many challenges

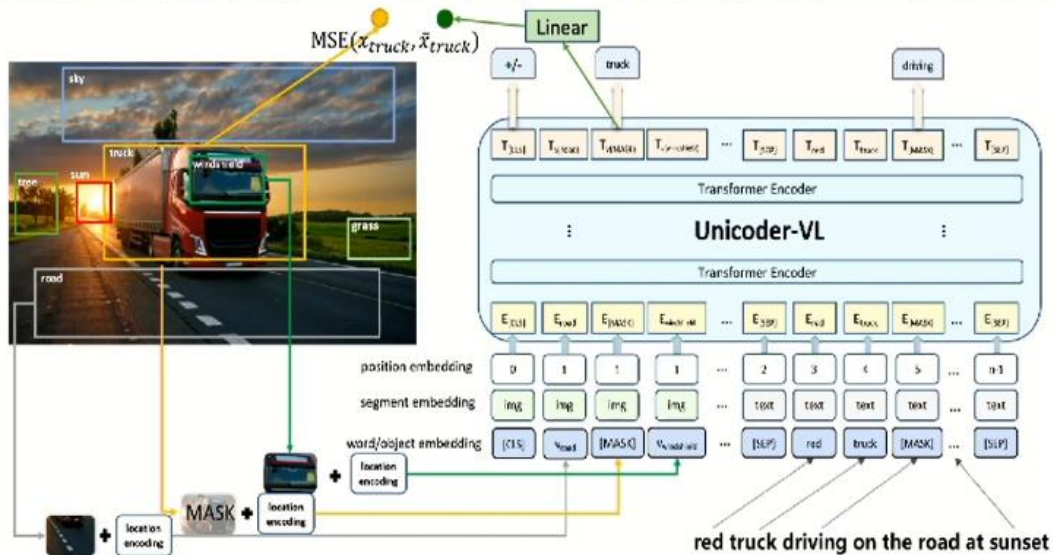
- Lacks training tasks fully leveraging the characteristics of multi-lingual data
- High cost to train huge parameters due to large vocab (250K)
- Transfer learning is weaker for languages in different families of language

# Agenda

1. Pre-trained models
2. Pre-trained models in multi-lingual tasks
3. Pre-trained models in multi-modality tasks
4. Summary

## Image-Language Pre-training

ViLBERT (Lu et al., 2019); Unicoder-VL (Li et al., 2019); VL-BERT (Su et al., 2019); UNITER (Chen et al., 2019)





# Evaluation Results

Model	Text-to-Image Retrieval (Flickr30k)			Image-to-Text Retrieval (Flickr30k)		
	R@1	R@5	R@10	R@1	R@5	R@10
VILBERT (Lu et al., 2019)	58.2	84.9	91.5	-	-	-
UNITER (Chen et al., 2019)	71.5	91.2	95.2	84.7	97.1	99.0
Unicoder-VL (Li et al., 2020)	73.1	92.3	95.9	88.0	97.3	96.6

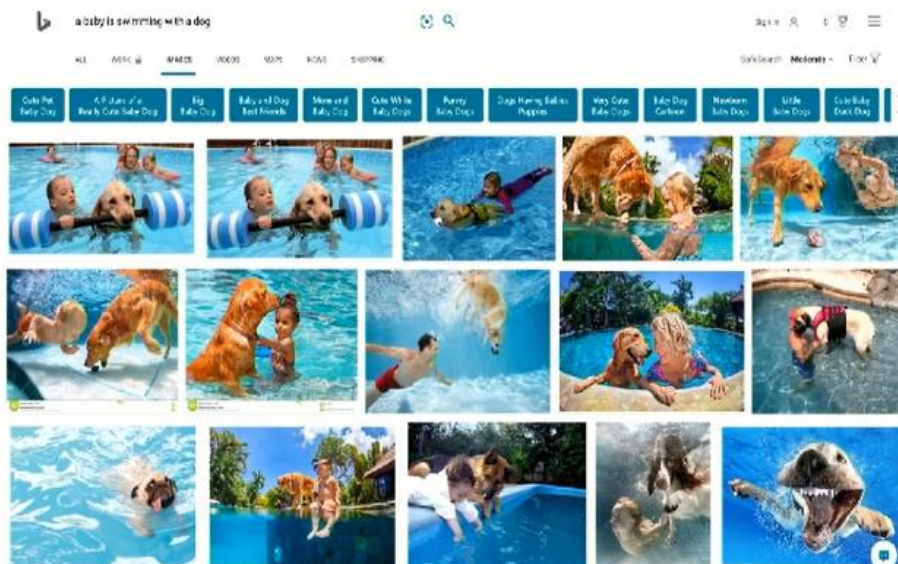
  

Model	Text-to-Image Retrieval (MSCOCO)			Image-to-Text Retrieval (MSCOCO)		
	R@1	R@5	R@10	R@1	R@5	R@10
UNITER (Chen et al., 2019)	48.4	76.7	85.9	63.3	87.0	93.1
Unicoder-VL (Li et al., 2020)	50.5	78.7	87.1	66.4	89.8	94.4

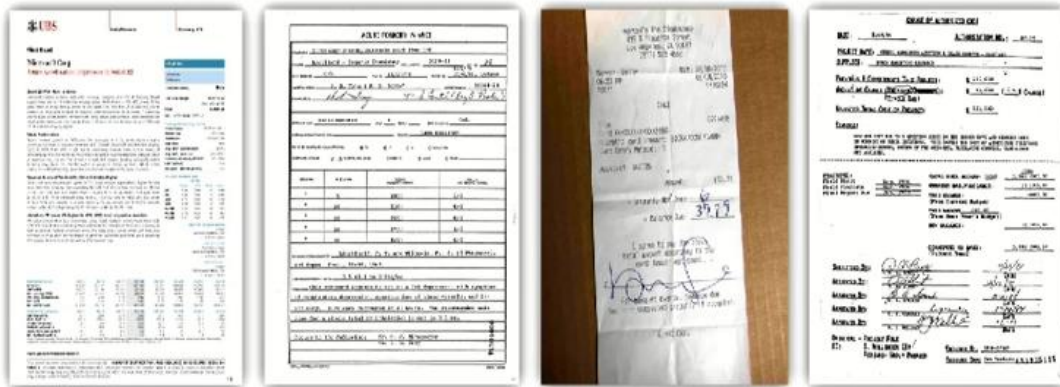
- Pre-training dataset
  - 3,318,333 image-caption pairs from Google's Conceptual Captions



# Application: Image Search



# Document Understanding



Information extraction from digital-born/scanned documents

LayoutLM: Pre-training for text with rich **Layout** and **Style** information  
<https://arxiv.org/pdf/1912.13318.pdf>, KDD 2020



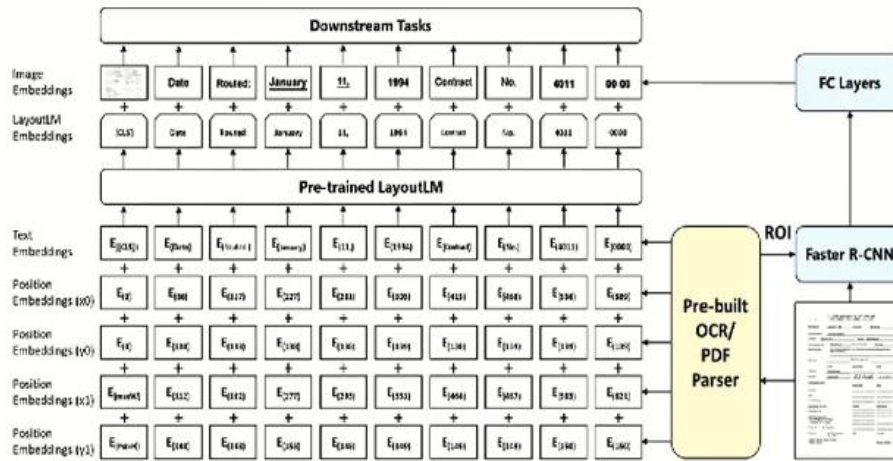
# Training Dataset for LayoutLM



11 million scanned document images from IIT-CDIP Test Collection 1.0  
<https://ir.nist.gov/cdip/>



# LayoutLM: Text+Layout Pre-training



<https://arxiv.org/abs/1912.13318>, KDD 2020

**SPORTS MARKETING ENTERPRISES  
DOCUMENT CLEARANCE SHEET**

Date Routed: January 11, 1994 Contract No. 4011 00 00

Contract Subject: Jon's Place Exhibits

Company: SPEVOO, INC. (Brand) Client/Version

Total Contract Cost: \$1,340,000.00 Contract Year: 1994-1998

Brief Description: 2 year's Place Exhibits for use at Weston Cup, Weston Dog and Crawl Boat Race Events.

GL Code: \_\_\_\_\_ Program Budget Code \_\_\_\_\_

NAME: \_\_\_\_\_ SIGNATURE: \_\_\_\_\_ DATE: \_\_\_\_\_

Original: Michael Wilkins

Message: John Powell [Signature] 1-11-94

RENEWAL/DUTY: \_\_\_\_\_ SIGNATURE: \_\_\_\_\_ DATE: \_\_\_\_\_

Insurance: \_\_\_\_\_

Law: \_\_\_\_\_

FBI - Marketing: \_\_\_\_\_

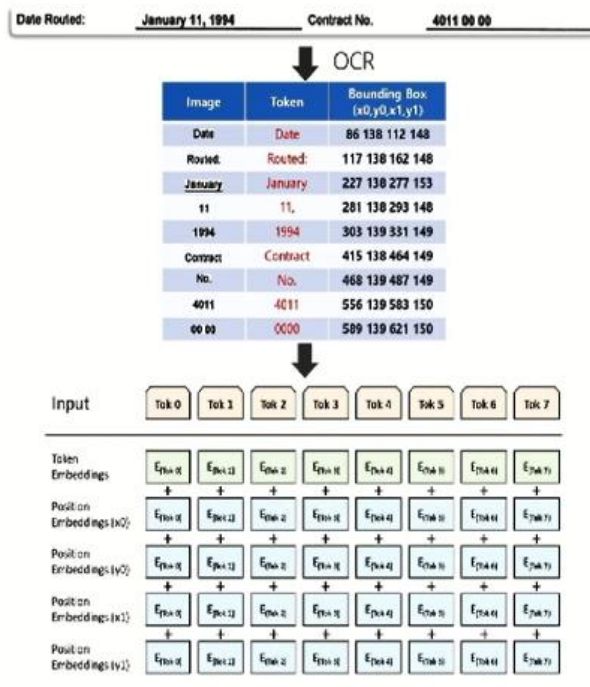
BEHAVIOR TO SHELL: \_\_\_\_\_ FACILITY: \_\_\_\_\_ SECTION: \_\_\_\_\_

APPROVAL SIGNATURES:  
 \* Sr. Manager (B. J. Powell)  
 \* Director - (G. L. Lohr)  
 \*\* Sr. VP - T. W. Robertson

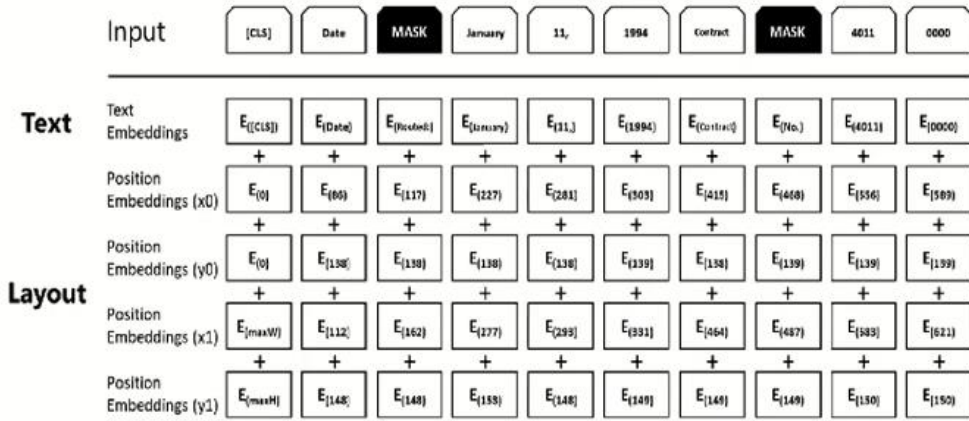
Revised To: MARY OSAGARVED DATE: 1/14/95 BY: [Signature]

\* UP TO AND INCLUDING \$25,000  
 \*\* OVER \$25,000

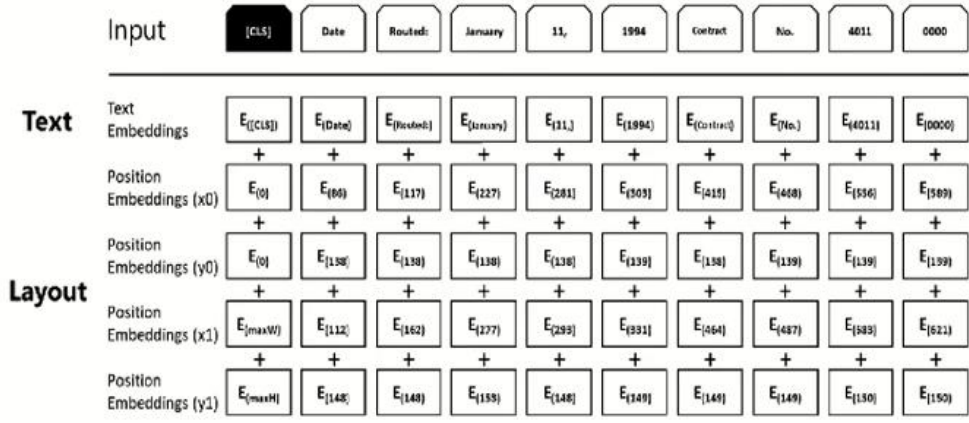
Revised 10/20/02



# Training Task#1: Mask Visual-Language Model



# Training Task#2: Document Image Classification



# Form Understanding with LayoutLM



**[Task]** Sequence labeling (B-I-O class labels) for key-value from forms  
**[Data]** 149 training, 50 testing  
**[Metric]** Precision, Recall, F1  
**[Baseline]** Pre-trained BERT and RoBERTa



FUNSD: Form Understanding in Noisy Scanned Documents  
<https://guillaumejaume.github.io/FUNSD/>

# Results on Form Understanding

Modality	Model	Precision	Recall	F1	Parameters
Text only	BERT <sub>BASE</sub>	0.547	0.671	0.603	110M
	RoBERTa <sub>BASE</sub>	0.632	0.69	0.66	125M
	BERT <sub>LARGE</sub>	0.611	0.709	0.656	340M
	RoBERTa <sub>LARGE</sub>	0.674	0.738	0.704	355M
[Base]Text+Layout	LayoutLM <sub>BASE</sub> (1M)	0.691	0.774	0.73	113M
	LayoutLM <sub>BASE</sub> (11M)	0.76	0.816	0.787	113M
[Base]Text+Layout+Image	LayoutLM <sub>BASE</sub> (1M)	0.71	0.782	0.744	160M
	LayoutLM <sub>BASE</sub> (11M)	0.768	0.82	0.793	160M
[Large]Text+Layout	LayoutLM <sub>LARGE</sub> (1M)	0.717	0.805	0.759	343M
	<b>LayoutLM<sub>LARGE</sub> (11M)</b>	<b>0.793</b>	<b>0.854</b>	<b>0.822</b>	343M
[Large]Text+Layout+Image	On-going	-	-	-	-

Pre-training data: IIT-CDIP Test Collection 1.0 <https://ir.nist.gov/cdip/>  
 Fine-tuning data: FUNSD dataset with 149 forms, test on 50 forms

## Results on Receipt Understanding

Modality	Model	Precision	Recall	F1	Parameters
Text only	BERT <sub>BASE</sub>	0.9099	0.9099	0.9099	110M
	RoBERTa <sub>BASE</sub>	0.9107	0.9107	0.9107	125M
	BERT <sub>LARGE</sub>	0.9200	0.9200	0.9200	340M
	RoBERTa <sub>LARGE</sub>	0.9280	0.9280	0.9280	355M
[Base]Text+Layout	LayoutLM <sub>BASE</sub> (1M)	0.9380	0.9380	0.9380	113M
	LayoutLM <sub>BASE</sub> (11M)	0.9438	0.9438	0.9438	113M
[Base]Text+Layout+Image	LayoutLM <sub>BASE</sub> (1M)	0.9416	0.9416	0.9416	160M
	LayoutLM <sub>BASE</sub> (11M)	0.9467	0.9467	0.9467	160M
[Large]Text+Layout	LayoutLM <sub>LARGE</sub> (1M)	0.9416	0.9416	0.9416	343M
	<b>LayoutLM<sub>LARGE</sub> (11M)</b>	<b>0.9604</b>	<b>0.9604</b>	<b>0.9604</b>	343M
Baseline	Ranking 1 <sup>st</sup> in ICDAR 2019	0.9402	0.9402	0.9402	-

ICDAR 2019 Robust Reading Challenge on Key Information Extraction from Scanned Receipts  
626 training, 347 testing, <https://rrc.cvc.uab.es/?ch=13&com=tasks>

## Results on Document Image Classification

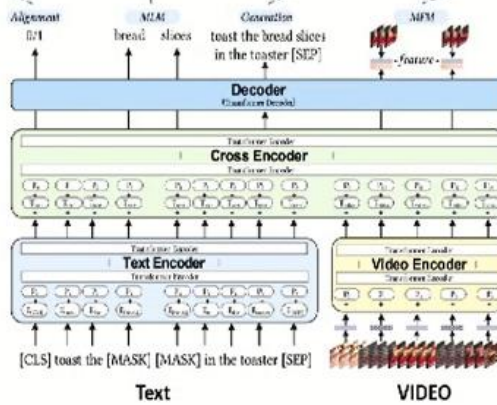
Modality	Model	ACC	Parameters
Text only	BERT <sub>BASE</sub>	89.81	110M
	RoBERTa <sub>BASE</sub>	90.06	125M
	BERT <sub>LARGE</sub>	89.92	340M
	RoBERTa <sub>LARGE</sub>	90.11	355M
[Base]Text+Layout	LayoutLM <sub>BASE</sub> (1M)	91.48	113M
	LayoutLM <sub>BASE</sub> (11M)	91.78	113M
[Base]Text+Layout+Image	LayoutLM <sub>BASE</sub> (1M)	94.31	160M
	<b>LayoutLM<sub>BASE</sub> (11M)</b>	<b>94.42</b>	160M
[Large]Text+Layout	LayoutLM <sub>LARGE</sub> (1M)	91.88	343M
	LayoutLM <sub>LARGE</sub> (11M)	91.90	343M
Previous SOTA	LadderNet (Sarkhel & Nandi, 2019)	92.77	-
	EfficientNetB4 (DeGange, 2019)	92.81	-
	MultiModal (Dauphinee et al., 2019)	93.07	-

RVL-CDIP dataset (320K training, 40K validation, 40K testing), <https://www.cs.cmu.edu/~aharley/rvl-cdip/>

# Video-Language Pre-training

VideoBERT (Sun et al., 2019); CBT (Sun et al., 2019); Unicoder-VL (Luo et al., 2020)

1. Video-Text Alignment 2. Masked Language Model 3. Transcript Generation 4. Masked Frame Model



# Video-Language Datasets

**Pre-training Data:**

1. HowTo100M: 136M video clips with captions from 1.2M Youtube videos.

**(Research) Evaluation Data:**

1. YouCook2: 2000 long untrimmed videos from 89 cooking recipes.
2. MSR-VTT: 10K web video clips with 41.2 hours and 200K clip-sentence pairs

**What is HowTo100M ?**

HowTo100M is a large-scale dataset of narrated videos with an emphasis on instructional videos where content creators teach complex tasks with an explicit intention of explaining the visual content on screen. HowTo100M features a total of:

- 136M video clips with captions sourced from 1.2M Youtube videos (15 years of video)
- 23k activities from domains such as cooking, hand-crafting, personal care, gardening or fitness.

Each video is associated with a narration available as subtitles automatically downloaded from Youtube.

<https://www.di.ens.fr/willow/research/howto100m/>

**YouCook2 Dataset** Home Explore Download Leaderboard

**Overview**

YouCook2 is the largest task-oriented, instructional video dataset in the open community. It contains 2000 long web videos from 89 cooking recipes, on average each of 10 minutes (vs. 25 videos). The procedure steps for each video are associated with temporal boundaries and described by imperative English sentences (see the example below). The videos were downloaded from Youtube and are all with a third-person viewpoint. All the videos are re-rendered and can be performed by individual persons at their homes with no need of camera. YouCook2 contains full video steps and videos cooking activities all over the world. It is the most complete and detailed dataset.

YouCook2 is currently suitable for video-language research, visual question-answering, and object recognition in videos, camera object and action discovery across videos and procedure learning.

**Note:** The release of the dataset includes pre-annotation for objects in the dataset. You can read more facts in <http://youcook2.eecs.umich.edu/>



# Evaluation Results

Unicoder-VL beats results from Google (CBT) and Baidu (ActBERT) on video retrieval and video captioning respectively.

Task	Video Retrieval						Video Captioning					
	YouCook2			MSR-VTT			YouCook2 (w/o transcript)			YouCook2 (w/ transcript)		
	R@1	R@5	R@10	R@1	R@5	R@10	BLEU-4	R-L	CIDEr	BLEU-4	R-L	CIDEr
SoTA	9.6	26.7	38.0	8.6	23.4	33.1	5.4	30.6	0.6	9.0	36.7	1.1
Unicoder-VL	10.0	27.5	38.8	15.4	39.5	52.3	6.1	31.5	0.6	10.4	38.0	1.2
$\Delta$	0.4 $\uparrow$	0.8 $\uparrow$	0.8 $\uparrow$	6.8 $\uparrow$	16.1 $\uparrow$	19.2 $\uparrow$	0.7 $\uparrow$	0.9 $\uparrow$	-	1.4 $\uparrow$	1.3 $\uparrow$	0.1 $\uparrow$

Blue numbers indicates the best result for a task.

### Clips



### Groundtruth

mince cabbage and chop some green onions

add soy sauce sesame oil and salt and mix together

### Our results

slice the cabbage green onions and add to the pot.

mix the soy sauce sesame oil and salt

# Application: Video Chaptering







相关阅读：

CCAI 2020 大会特邀报告

CCAI 2020 大会特邀报告(II): Thoughtful Artificial Intelligence



临菲信息技术港



临菲信息技术港公众号



临菲学堂