



# 大语言模型

Large Language Model (LLM)





随着近年来人工智能技术的飞速发展,大模型的发展趋势也可谓是"激流勇进",但在迅猛发展的背后,技术的局限性也日益凸显。研究者们针对这些局限性也在积极探索并提出了不计其数的新方法和研究方向。

# 章节概述

CHAPTER OVERVIEW



**超大规模预训练网络的基本原理及典型网络代表** 

**主流模型压缩方法** 



**其他自然语言处理热门研究点** 



问答系统

机器阅读理解

→ 多模态任务

# 小节介绍

SECTION INTRODUCTION



# 超大规模预训练网络

预训练技术, BERT, GPT系列

2

模型压缩方法

模型剪枝,模型量化,模型蒸馏

3

其它热门的研究点

问答系统, 机器阅读理解

4

多模态任务的举例与现状

多模态学习,图像-文本多模态任务

预训练是通过设计好一个网络结构来做语言模型任务,然后把大量甚至是无穷尽的无标注自然语言文本利用起来,预训练任务把大量语言学知识抽取出来编码到网络结构中,在当前任务带有标注信息的数据有限时,这些先验的语言学特征会对当前任务有极大的特征补充作用。



# 大语言模型(Large Language Model, LLM)发展史





**GPT** 06/2018



GPT-2 02/2019



GPT-3 05/2020



**GPT-4** 05/2023

#### Transformer 06/2017

#### **BERT** 10/2018





### ●大语言模型的"三大"特点:

- ●大数据
- ●大模型
- ●大算力



天河天元大模型

# 大语言模型(Large Language Model, LLM)

# 国内外大模型百花齐放



◆ 360智脑

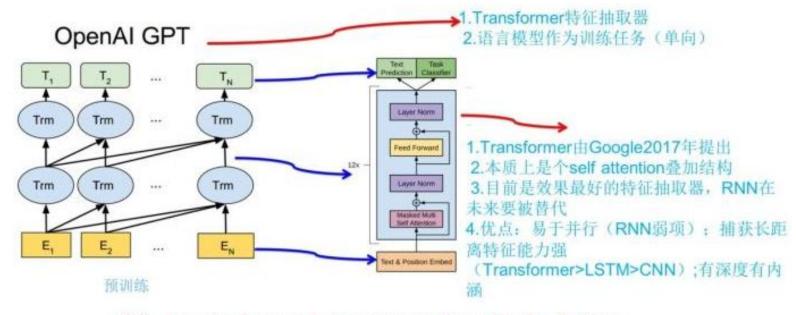
★ METASOTA A 写作猫





# **Generative Pre-training Transformer (GPT)**

从WE到GPT: Pretrain+Finetune两阶段过程



"基于微调的模式

论文: Improving Language Understanding by Generative Pre-Training



第一阶段 → 利用语言模型进行预训练

通过微调的模式解决下游任务



# Generative Pre-Training (GPT) 预训练方式

Transformer在做机器翻译任务时,需要平行数据集。但是我们身边存在大量没有标注的数据,例如文 本、图片、代码等等。标注这些数据需要花费大量的人力和时间,标注的速度远远不及数据产生的速度, 所以带有标签的数据往往只占有总数据集很小的一部分。随着算力的不断提高,计算机能够处理的数据量 逐渐增大。如果不能很好利用这些无标签的数据就显得很浪费。

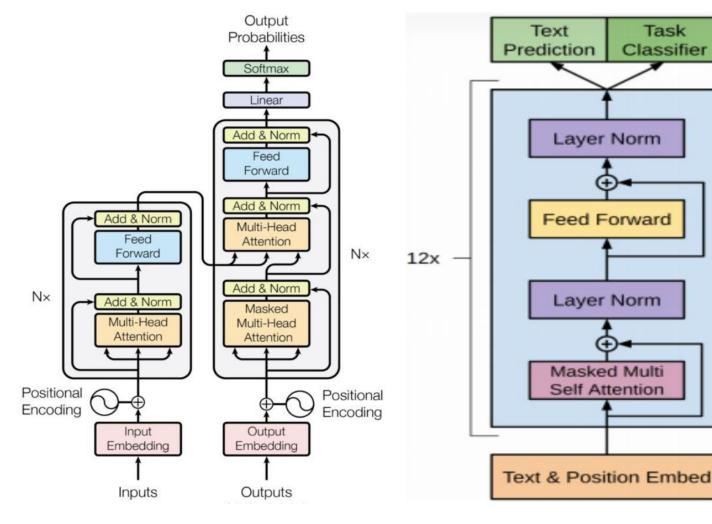
半监督学习和预训练+微调的二阶段模式整变得越来越受欢迎。最常见的二阶段方法就是Word2Vec,使用 大量无标记的文本训练出带有一定语义信息的词向量,然后将这些词向量作为下游机器学习任务的输入, 就能够大大提高下游模型的泛化能力。但是Word2Vec有一个问题,就是单个单词只能有一个Embedding。 这样一来,一词多义就不能很好地进行表示。



# Generative Pre-Training (GPT) 预训练方式

# 单向Transformer结构

在Transformer的文章中,提到了 Encoder与Decoder使用的 Transformer Block是不同的。在 Decoder Block中,使用了Masked Self-Attention,即句子中的每个词,都 只能对包括自己在内的前面所有词进行 Attention, 这就是单向Transformer。 GPT使用的Transformer结构就是将 Encoder中的Self-Attention替换成了 Masked Self-Attention.







Task

Classifier

Layer Norm

Feed Forward

Laver Norm

Masked Multi Self Attention





# Generative Pre-Training (GPT) 训练方式

### GPT预训练过程

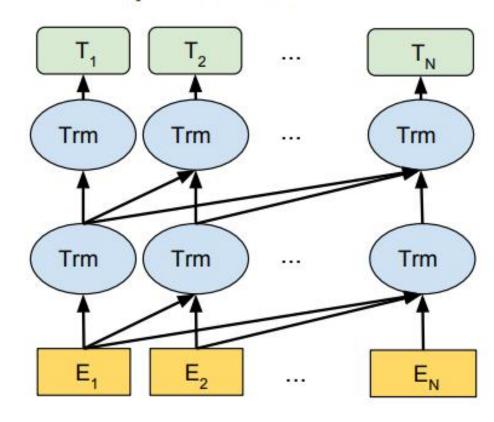
由于采用的是单向的Transformer,只能看到上文的词, 所以语言模型为:

$$L_1(\mathcal{U}) = \sum_{i} \log P(u_i|u_{i-k}, \dots, u_{i-1}; \Theta)$$

训练过程就是将句子n个词的词向量(第一个为<BOS>)加上 Positional Encoding后输入到前面提到的Transfromer中, n个输出分别预测该位置的下一个词(<BOS>预测句子 中的第一个词, 最后一个词的预测结果不用于语言模型的 训练)。

由于使用了Masked Self-Attention, 所以每个位置的词 都不会"看见"后面的词,也就是预测的时候是看不见 "答案"的,保证了模型的合理性,这也是为什么OpenAl 采用了单向Transformer的原因。

# OpenAl GPT





# **Generative Pre-Training (GPT) 微调过程**

接下来就进入模型训练的第二步,运用少量的带标签数据对模型参数进行微调。

上一步中最后一个词的输出我们没有用到,在这一步中就要使用这一个输出来作为下游监督学习的 输入。

$$L_2(\mathcal{C}) = \sum_{(x,y)} \log P(y|x^1,\dots,x^m).$$

$$P(y|x^1,\ldots,x^m) = \operatorname{softmax}(h_l^m W_y).$$

为避免Fine-Tuning使得模型陷入过拟合,文中还提到了辅助训练目标的方法,类似于一个多任务 模型或者半监督学习。具体方法就是在使用最后一个词的预测结果进行监督学习的同时,前面的词 继续上一步的无监督训练, 使得最终的损失函数成为:

$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C})$$



# 大语言模型(Large Language Model, LLM)发展史





**GPT** 06/2018



GPT-2 02/2019



GPT-3 05/2020



**GPT-4** 05/2023

#### Transformer 06/2017

#### **BERT** 10/2018





- ●大语言模型的"三大"特点:
  - ●大数据
  - ●大模型
  - ●大算力





#### BERT与GPT的相似之处:

BERT(Bidirectional Encoder Representation from Transformer), 是Google Brain在2018年提 出的基于Transformer的自然语言表示框架。是一提出就大火的明星模型。BERT与GPT一样,采取了 Pre-training + Fine-tuning的训练方式,在分类、标注等任务下都获得了更好的效果。 BERT与GPT非常的相似,都是基于Transformer的二阶段训练模型,都分为Pre-Training与Fine-Tuning两个阶段,都在Pre-Training阶段无监督地训练出一个可通用的Transformer模型,然后在 Fine-Tuning阶段对这个模型中的参数进行微调,使之能够适应不同的下游任务。

#### BERT与GPT的不同之处:

训练目标和模型结构和使用上不同:

GPT采用的是单向的Transformer, 而BERT采用的是双向的Transformer, 也就是不用进行Mask操作; 使用的结构的不同,直接导致了它们在Pre-Training阶段训练目标的不同;

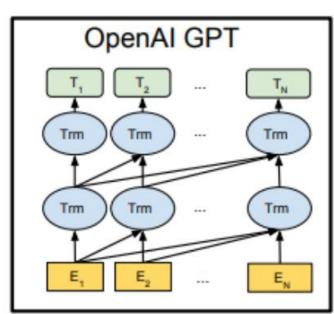




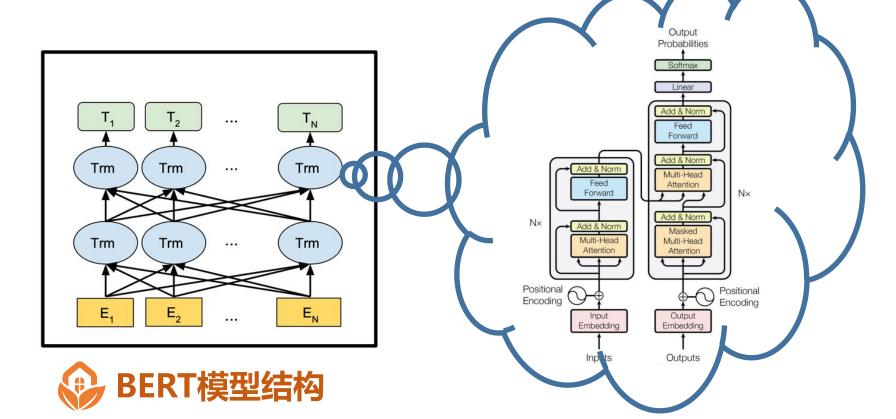
# 基于Transformer的双向编码器表示

预训练方法的创新

用了Masked LM和Next Sentence Prediction 两种方法分别捕捉词语和句子级别的表征











## Masked Language Model (MLM)

在Transformer中,我们即想要知道上文的信息,又想要知道下文的信息,但同时要保证整个模型不知道要预测词的信息,那么就干脆不要告诉模型这个词的信息就可以了。也就是说,BERT在输入的句子中,挖掉一些需要预测的词,然后通过上下文来分析句子,最终使用其相应位置的输出来预测被挖掉的词。这其实就像是在做完形填空 (Cloze)一样。

但是,直接将大量的词替换为<MASK>标签可能会造成一些问题,模型可能会认为只需要预测 <MASK>相应的输出就行,其他位置的输出就无所谓。同时Fine-Tuning阶段的输入数据中并没有 <MASK>标签,也有数据分布不同的问题。为了减轻这样训练带来的影响,BERT采用了如下的方式:

- 1. 输入数据中随机选择15%的词用于预测,这15%的词中,
- 2.80%的词向量输入时被替换为<MASK>
- 3. 10%的词的词向量在输入时被替换为其他词的词向量
- 4. 另外10%保持不动

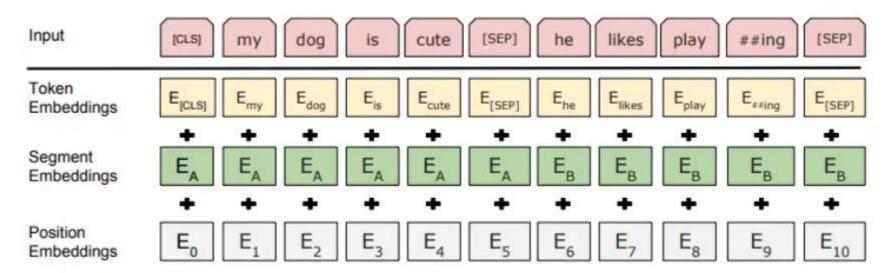
这样一来就相当于告诉模型,我可能给你答案,也可能不给你答案,也可能给你错误的答案,有 <MASK>的地方我会检查你的答案,没<MASK>的地方我也可能检查你的答案,所以<MASK> 标签对你来说没有什么特殊意义,所以无论如何,你都要好好预测所有位置的输出。



#### **Next Sentence Prediction (NSP)**

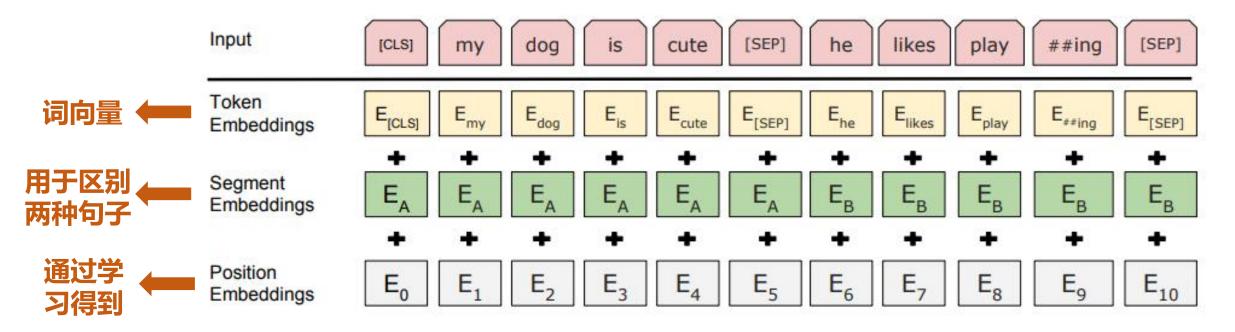
BERT还提出了另外一种预训练方式NSP,与MLM同时进行,组成多任务预训练。这种预训练的方式就是往Transformer中输入连续的两个句子,左边的句子前面加上一个<CLS>标签,它的输出被用来判断两个句子之间是否是连续上下文关系。采用负采样的方法,正负样本各占50%。

为了区分两个句子的前后关系,BERT除了加入了Positional Encoding之外,还两外加入了一个在预训练时需要学习的Segment Embedding来区分两个句子。这样一来,BERT的输入就由词向量、位置向量、段向量三个部分相加组成。此外,两个句子之间使用<SEP>标签予以区分。





# **BERT的Embedding组成**



# 第一步预训练 一

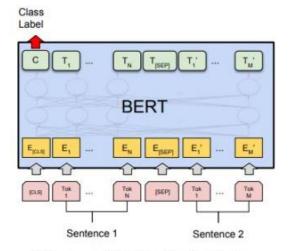
加mask的技巧 → 训练过程中随机mask 15%的token → 最终的损失函数只计算被mask掉的那个token

第二步预训练 — 让模型理解两个句子之间的联系

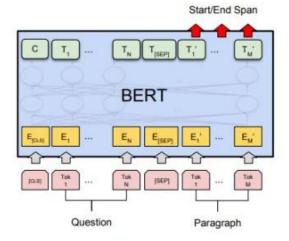


# **其他任务的BERT参数调整**

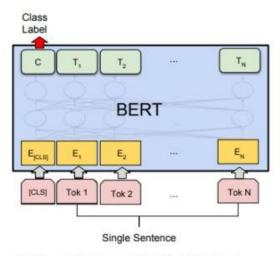
# Fine-tunning阶段



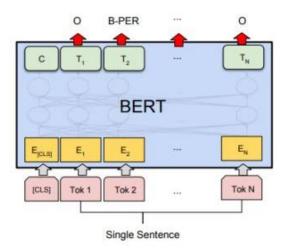
(a) Sentence Pair Classification Tasks: MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG



(c) Question Answering Tasks: SQuAD v1.1



(b) Single Sentence Classification Tasks: SST-2, CoLA



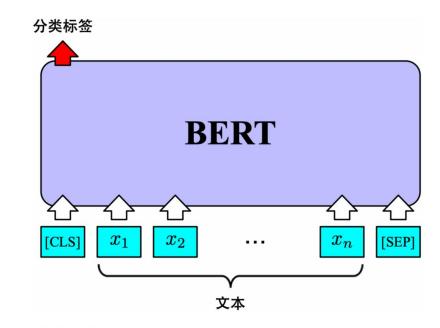
(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

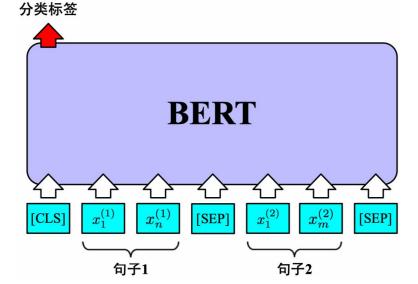


# 下游任务应用

- ●预训练语言模型应用
  - 单句文本分类
    - 最常见的NLP任务,需要将输入文本分成不同类别
    - 例如: 情感分类等

- ●句对文本分类
  - 与单句分类任务类似,需要将一对文本分成不同类别
  - 例如: 文本蕴涵任务中,将句对分成"蕴涵"或"冲突" 类别







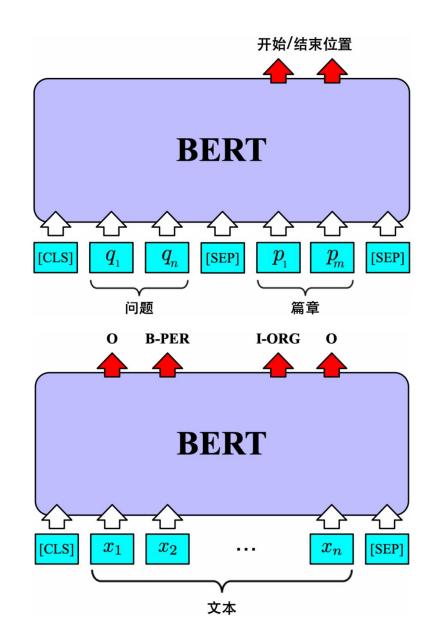
#### 下游任务应用

- ●预训练语言模型应用
  - ●阅读理解
    - 以抽取式阅读理解为例,要求机器在阅读篇章和问题后 给出相应的答案, 而答案要求是从篇章中抽取出一个文 本片段(Span)

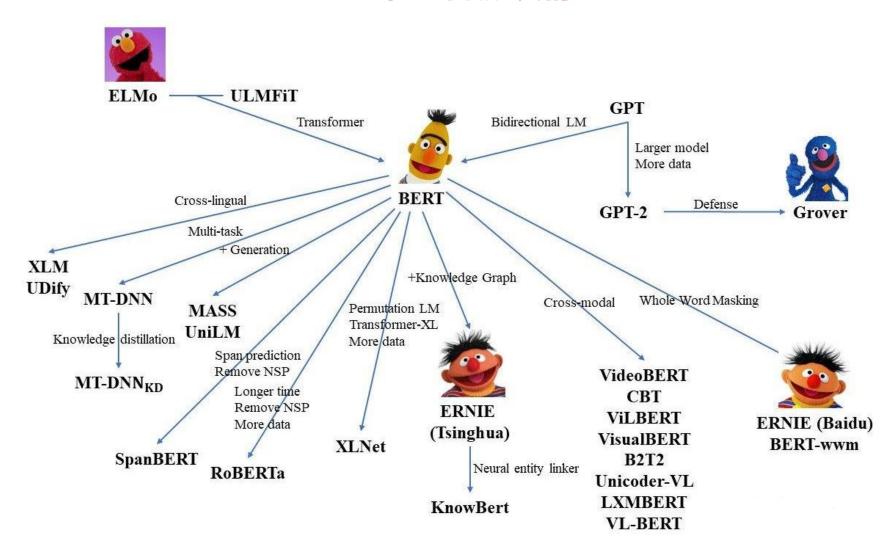
#### ● 序列标注

● 以命名实体识别任务为例,对给定输入文本的每一个词输出 一个标签, 以此指定某个命名实体的边界信息

原始标签	B-PER	I-PER	О	О	О	О	B-LOC	
原始输入	John	Smith	has	never	been	to	Harbin	
处理后的标签	B-PER	I-PER	О	О	О	О	B-LOC	B-LOC
<u></u>	John	Smith	has	never	been	to	Ha	##rbin



### BERT和它的朋友们





#### **RoBERTa**

- ●在BERT的基础上引入了动态掩码技术,同时舍弃了NSP任务
- ●采用了更大规模的预训练数据,并以更大的批次和BPE词表训练了更多的步数
- ●动态掩码
  - 即决定掩码位置和方法是在模型的训练阶段实时计算的,这样能保证无论训练多少轮,都 能够最大限度地保证同一段文本能够在不同轮数下产生不同的掩码模型

```
went to the store → went to the [MASK]
went to the store → went to [MASK] store
went to the store → go to the store
went to the store → went to the store
went to the store → went to the store
went to the store → went [MASK] the store

动态掩码
```

#### **RoBERTa**

- ●舍弃NSP任务
  - 实验结果表明不使用NSP任务能够带来下游任务的性能提升
  - ●最后采用了"跨文档整句输入"并舍弃了NSP任务的方案

实验设置	SQuAD 1.1	SQuAD 2.0	MNLI-m	SST-2	RACE
文本对输入 + NSP	90.4	78.7	84.0	92.9	64.2
句子对输入 + NSP	88.7	76.2	82.9	92.1	63.0
跨文档整句输入	90.4	79.1	84.7	92.5	64.8
文档内整句输入	90.6	79.7	84.7	92.7	65.6

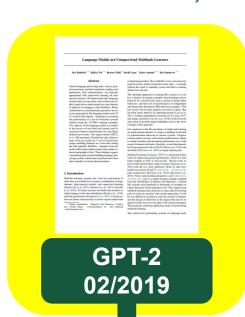


# 大语言模型(Large Language Model, LLM)发展史





**GPT** 06/2018







#### Transformer 06/2017

#### **BERT** 10/2018





### ●大语言模型的"三大"特点:

- ●大数据
- ●大模型
- ●大算力





GPT-2继续沿用了原来在GPT种使用的单向Transformer模型,目的就是尽 可能利用单向Transformer的优势,做一些BERT使用的双向Transformer所 做不到的事。那就是通过上文生成下文文本。

GPT-2的想法就是完全舍弃Fine-Tuning过程,转而使用一个容量更大、无监督 训练、更加通用的语言模型来完成各种各样的任务。我们完全不需要去定义这 个模型应该做什么任务,因为很多标签所蕴含的信息,就存在于语料当中。就 像一个人如果博览群书,自然可以根据看过的内容轻松的做到自动摘要、问答、 续写文章这些事。

严格来说GPT-2可能不算是一个多任务模型,但是它确实使用相同的模型、相 同的参数完成了不同的任务。那么GPT-2是怎么使用语言模型完成多种任务的 呢?



# 主流超大规模预训练网络介绍-GPT-2



通常我们针对特定任务训练的专用模型,给定输入,就可以返回这个任务相应的输出,也就是

# $p(output \mid input)$

那么如果我们希望设计一个通用的模型,这个模型在给定输入的同时还需要给定任务类型,然后根 据给定输入与任务来做出相应的输出,那么模型就可以表示成下面这个样子

### $p(output \mid input, task)$

就好像原来我需要翻译一个句子,需要专门设计一个翻译模型,想要问答系统需要专门设计一个问题。 答模型。但是如果一个模型足够聪明,并且能够根据你的上文生成下文,那我们就可以通过在输入 中加入一些标识符就可以区分各种问题。比如可以直接问他: ('自然语言处理',中文翻译)来得到 我们需要的结果Nature Language Processing。在我的理解中GPT-2更像是一个无所不知的问答 系统,通过告知一个给定任务的标识符,就可以对多种领域的问答、多种任务做出合适的回答。 GPT-2满足零样本设置 (zero-shot setting), 在训练的过程中不需要告诉他应该完成什么样的任 务, 预测是也能给出较为合理的回答。



性能稳定、优异。

有着超大规模的在海量数据集上训练的基于Transformer的巨大模型 模型与只带有解码器的Transformer模型类似



GPT-2 GPT-2将Transformer堆叠的层数增加到48层,隐层的维度为1600, 参数量达到了15亿。

> GPT-2将词汇表提升到50257。 处理最长单词序列数由GPT的512扩展到1024。 此外还对Transformer做出了小调整。



# 主流超大规模预训练网络介绍-GPT-2



# GPT-2中Transformer模块的处理方式

每个Transformer模块的处理方式都是一样的,但每个模块都会维护自己的自注意 力层和神经网络层中的权重

通过自注意力层处理



将其传递给神经网络层



模块处理完单词后,将结果向量传入堆栈中的下 一个Transformer模块,继续进行计算



最后一个Transformer模块产生输出后,模型会将输出的向 量乘上嵌入矩阵得到词汇表中每个单词对应的注意力得分





GPT-2的最大贡献是验证了通过海量数据和大量参数训练出来的词向量模型有迁移到其它类别任务中而不需要额外的训练。但是很多实验也表明,GPT-2的无监督学习的能力还有很大的提升空间,甚至在有些任务上的表现不比随机的好。尽管在有些zero-shot的任务上的表现不错,但是我们仍不清楚GPT-2的这种策略究竟能做成什么样子。GPT-2表明随着模型容量和数据量的增大,其潜能还有进一步开发的空间,基于这个思想,诞生了我们下面要介绍的GPT-3。



# 大语言模型(Large Language Model, LLM)发展史

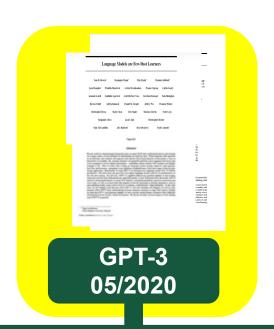




**GPT** 06/2018



GPT-2 02/2019





Transformer 06/2017

**BERT** 10/2018





### ●大语言模型的"三大"特点:

- ●大数据
- ●大模型
- ●大算力



# 主流超大规模预训练网络介绍-GPT-3



模型	发布时间	参数量	预训练数据量
GPT	2018年6月	1.17 亿	约 5GB
GPT-2	2019年2月	15 亿	40GB
GPT-3	2020年5月	1,750 亿	45TB

GPT-3沿用GPT-2的网络结构。为了使模型在少样本的任务中尽快收敛,提出In-context learning的概念,核心思想在于通过少量的数据寻找一个合适的初始化范围,使得模型能够在有限的数据集上快速拟合,并获得不错的效果。



# 主流超大规模预训练网络介绍



#### **OpenAI Internal Timeline By FeltSteam**





# 主流超大规模预训练网络介绍



GPT-6也被电力卡脖子了——部署十万个H100时,整个电网发生了崩溃!

就在刚刚,微软工程师爆料,10万个H100基建正在紧锣密鼓地建设中,目的就是训练GPT-6.

微软工程师吐槽说,团队在部署跨区域GPU间的infiniband级别链接时遇到了困难。

Corbitt: 为何不考虑直接将所有设备部署在同一个地区呢?

微软工程师:这确实是我们最初的方案。但问题是,一旦我们在同一个州部署超过100,000 个H100 GPU, 电网就会因无法负荷而崩溃。



1

#### 超大规模预训练网络

预训练技术, BERT, GPT-2

2

#### 模型压缩方法

模型剪枝,模型量化,模型蒸馏

3

#### 其它热门的研究点

问答系统, 机器阅读理解

4

#### 多模态任务的举例与现状

多模态学习,图像-文本多模态任务

复杂的模型同时带来了高额的存储空间及计算资源消耗,使其较难落实到各个硬件平台。为了解决这些问题,需对模型进行模型 压缩以最大限度地减小模型对计算空间和时间的消耗。







### ● 基于的假设 (共识) → DNN的过参数化 (Over-parameterization)

DNN的过参数化: 指训练阶段网络需要大量的参数来捕捉数据中的微小信 息,而当训练完成并进入预测阶段后,网络通常并不需要这么多的参数。



#### 基本思想



#### 剪裁最不重要的部分



贪心法

考虑参数裁剪对损失的影响

考虑对特征输出的可重建性的影响



② 贪心法 (saliency-based方法) → 按重要性进行排序,之后将不重要的部分去除

magnitude-based weight pruning方法 小来评估重要性,然后用贪心法

按参数 (或特征输出) 绝对值大 对重要性较低的部分进行剪枝

结构化剪枝 — 用Group LASSO算法来得到结构化的稀疏权重

缺点! >> 忽略了参数间的相互关系而只能找到局部最优解



考虑参数裁剪对损失的影响



参 考虑对特征输出的可重建性的影响 → 最小化裁剪后网络对于特征输出的重建误差

如果对当前层进行裁剪后对后面的输出没有较大影响,则 基于的认知 说明裁掉的是不太重要的信息





### **→** 量化模型 (Quantized Model) → 在硬件上移植非常方便

- → 模型加速 (Model Acceleration) 方法中其中一类方法的总称
- 包括二值化网络 (Binary Network) 、三值化网络 (Ternary Network),深度压缩 (Deep Compression)等
- **通往高速神经网络最佳的方法,但仍面临实现难度大、准确性** 不稳定,使用门槛较高的多方面问题



量化后的权值张量 —— 一个高度稀疏的有较多共享权值的矩阵

对于非零参数 - 定点压缩 - 更高的压缩率

e.g. Deepcompression





### ₩ 采用的方法 → 迁移学习



通过预先训练好的复杂模型(Teacher model)的最后输出

→ 结果来作为先验知识,结合One-Hot label数据,共同指导一 个简单的网络(Student model)学习



**让student学习到teacher的泛化能力** 



hard target

temperature

# 常用名词解释

原始模型或模型ensemble teacher

student 新模型

用来迁移teacher知识、训练student的数据集合 transfer set

teacher输出的预测结果 (一般是Softmax之后的概率) soft target

样本原本的标签

蒸馏目标函数中的超参数

蒸馏的一种,指student和teacher的结构和尺寸完全一样

防止student的表现被teacher限制,在蒸馏时逐渐减少 soft targets的权重

born-again network teacher annealing



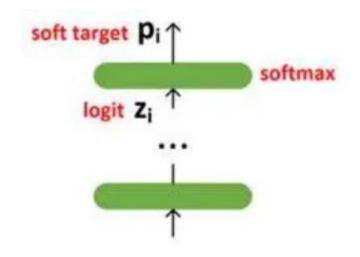


#### 原始模型训练阶段 ---

# 根据目标问题,设计一个大模型或者多个模型集合(N1, N2,...,Nt)即teacher,然后并行训练集合中的网络



#### 精简模型训练阶段



$$q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

 $z_i$ : Logit

例如: SoftMax层的输入

 $q_i$ : SoftMax层计算的分类概率

T: 温度, 通常被设为1





#### 精简模型训练阶段

设计一个简单网络NO即student并收集简单模型训练数据







对所有概率向量求取均值作为当前样本最后的概率输出向量并保存



标签融合前面收集到的数据定义为样本原本的标签,即hard\_target





1

#### 超大规模预训练网络

预训练技术, BERT, GPT-2

2

#### 模型压缩方法

模型剪枝,模型量化,模型蒸馏

3

#### 其它热门的研究点

问答系统, 机器阅读理解

4

#### 多模态任务的举例与现状

多模态学习,图像-文本多模态任务

自然语言处理的具体表现形式包括机器翻译、文本摘要、文本分类、文本校对、信息抽取、语音合成、语音识别等。自然语言处理就是要计算机理解自然语言,自然语言处理机制涉及两个流程,包括自然语言理解和自然语言生成。



信息检索系统的一种高级形式,通过Web搜索或链接知识

→ 问答系统 → 库等方式,检索到用户问题的答案,并用准确、简洁的自

然语言回答用户

分类-根据其问题答案的数据来源和回答的方式

基于Web信息检索的问答系统 (Web Question Answering, WebQA)

基于知识库的问答系统(Knowledge Based Question Answering, KBQA)

**社区问答系统(Community Question Answering, CQA)** 





#### 分类-根据其问题答案的数据来源和回答的方式

- **◇ 基于Web信息检索的问答系统 (Web Question Answering, WebQA)** 
  - 以搜索引擎为支撑,理解分析用户的问题意图后并在全网范围内搜索 相关答案反馈给用户
    - e.g. 早期的 Ask Jeeves 和 AnswerBus 问答系统
- 基于知识库的问答系统(Knowledge Based Question Answering, KBQA)
  - 通过结合一些已有的知识库或数据库资源,以及利用非结构化文本的信息,使用信息抽取的方法提取有价值的信息,并构建知识图谱作为问答系统的后台支撑,再结合知识推理等方法为用户提供更深层次语义理解的答案
- **社区问答系统(Community Question Answering, CQA)** 
  - 基于社交媒体的问答系统,大多数问题的答案由网友提供,问答系统会检索社交媒体中与用户提问语义相似的问题,并将答案返回给用户





- 知识库 (Knowledge Base, KB) 用于知识管理的一种特殊的数据库,用于相关领域知识的采集、整理及提取
  - 表示形式是一个对象模型 (object model) ,通常称为本体,包含一些类、子类和实体
  - 常见的知识库有Freebase、DBPedia等
  - 一般采用RDF格式对其中的知识进行表示,知识的查询主要采用RDF 标准查询语言SPARQL





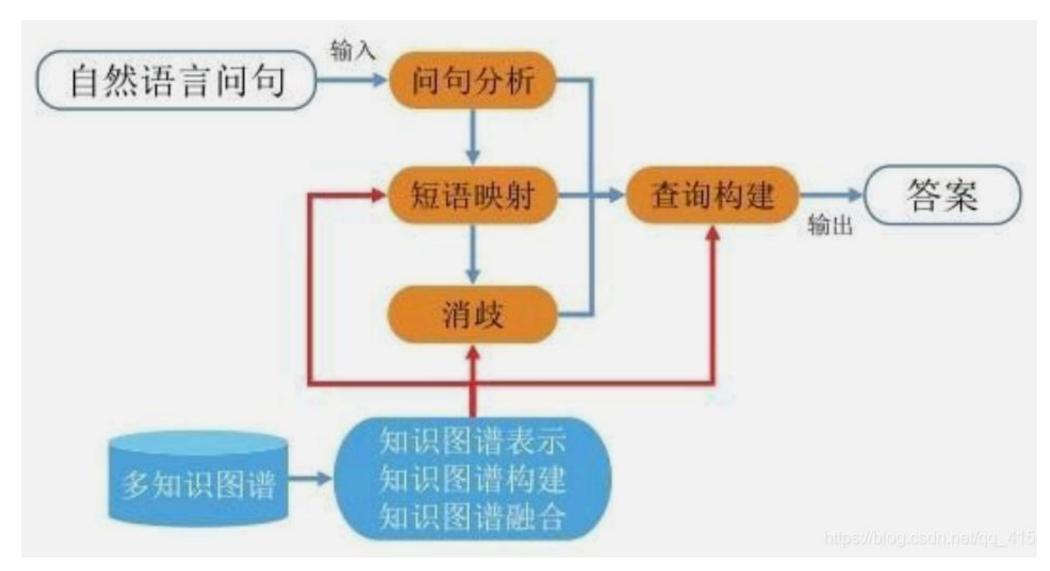
- 知识库 (Knowledge Base, KB) 用于知识管理的一种特殊的数据库,用于相关领域知识的采集、整理及提取
  - 表示形式是一个对象模型 (object model) ,通常称为本体,包含一些类、子类和实体
  - 常见的知识库有Freebase、DBPedia等
  - 一般采用RDF格式对其中的知识进行表示,知识的查询主要采用RDF 标准查询语言SPARQL



#### 热门研究点介绍-问答系统



#### KBQA基本架构 ➡️ 包含问句理解、答案信息抽取、答案排序和生成等核心模块





#### 热门研究点介绍-机器阅读理解



#### 机器阅读理解 (Machine Reading Comprehension, MRC)

**→** 给定一篇文章以及基于文章的一个问题,让机器在阅读文章后对问题进行作答



#### 常见任务主要







多项选择

给定一篇文章和一个问题,让模型从多个备选答案中选 择一个或多个最有可能是正确答案的选项



分 片段抽取

给定一篇文章和一个问题,让模型从文章中抽取连续的 单词序列,并使得该序列尽可能的作为该问题的答案



自由作答

给定一篇文章和一个问题,让模型生成一个单词序列, 并使得该序列尽可能的作为该问题的答案





#### **经典机器阅读理解的基本框架**



**特征抽取 (Feature Extraction、Encode)** 

文章-问题交互(Context-Question Interaction)

**答案预测 (Answer Prediction)** 



# 研究趋势

- 基于外部知识的机器阅读理解 → 相关外部知识的检索及外部知识的融合
- 不能回答的问题的判别、干扰答 带有不能回答的问题的机器阅读理解 —— 室的识别等
- 相关文档的检索、噪声文档的干扰、检索得到 多条文档机器阅读理解 → 的文档中没有答案、可能存在多个答案、需要 对多条线索进行聚合等
- 对话式阅读理解 → 对话历史信息的利用、指代消解等



1

#### 超大规模预训练网络

预训练技术,BERT, GPT-2

2

#### 模型压缩方法

模型剪枝,模型量化,模型蒸馏

3

#### 其它热门的研究点

问答系统, 机器阅读理解

4

### 多模态任务的举例与现状

多模态学习,图像-文本多模态任务

每一种信息的来源或者形式,都可以称为一种模态。例如,人有触觉,听觉,视觉等;信息的媒介,有语音、视频、文字等。以上的每一种都可以称为一种模态。多模态机器学习旨在通过机器学习的方法实现处理和理解多源模态信息的能力。目前比较热门的研究方向是图像、视频、音频、语义之间的多模态学习。

















多模态表示学习。

而学习到更好的特征表示

利用多模态之间的互补性,剔除模态间的冗余性,从

两大研究方向

联合表示 (Joint Representations)

协同表示 (Coordinated Representations)

○联合表示: 特征表示 多模态向量空间 模态 2 模态 1

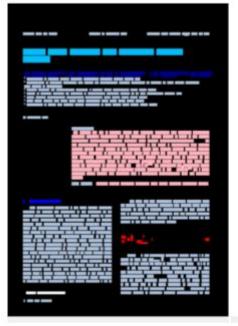
的相关性约束 协同表示 特征表示 2 特征表示 模态 1 模态 2

将多模态中的



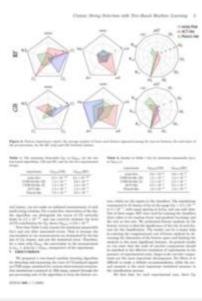
# **Document Al**

海量纸质数据需要电子化存档,通用的做法是先将纸质文件进行扫描,得到扫描文本图像,将文本图像进行OCR识别(计算机视觉,CV),得到文本信息,再对文本信息进行信息抽取和结构化存储(自然语言处理,NLP)。该任务在教育、金融、医疗等领域均有广泛的应用价值,因而得到了广泛的关注和研究,并形成了大数据、人工智能领域一个热点研究方向:文档智能(Document AI)。微软亚洲研究院在该领域提出了文档多模态预训练模型LayoutLM系列,得到了广泛的关注。利用多模态预训练模型LayoutLM可对多种文档领域下游任务带来提升,比如文档分类,版面分析,信息提取,文档VQA(视觉问答)等。



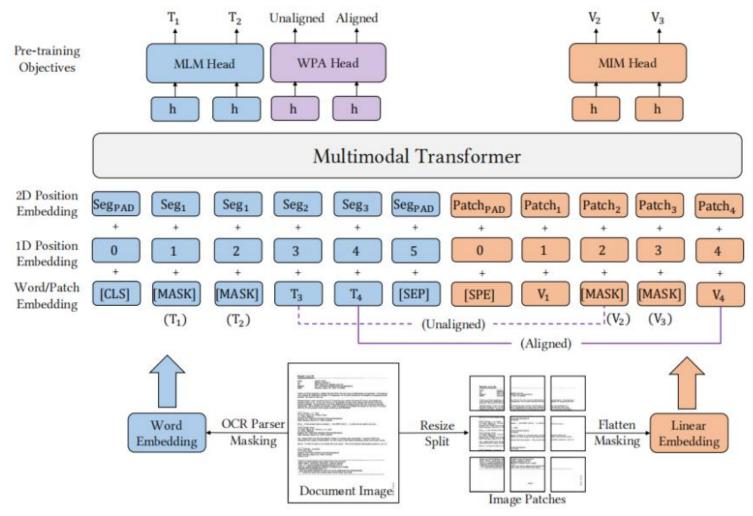








# **Document Al**



多模态模型LayoutLM,处理图像、文本和版式信息。





→ 模态转化(映射) → 将一个模态的信息转换为另一个模态的信息







多模态融合 ➡ 联合多个模态的信息,进行目标预测(分类或者回归)



分别对应对原始数据进行融合、对抽象的特征 进行融合和对决策结果进行融合

e.g. 迁移学习







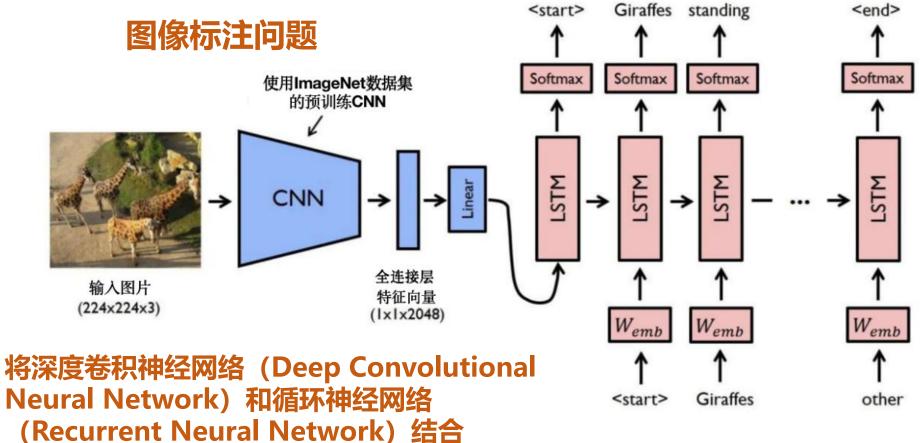
#### **❷** 图像描述 (Image Caption) ■

# 学习的综合问题

## 融合计算机视觉、自然语言处理和机器



#### 图像描述主体网络结构





#### 图像-文本多模态任务举例及研究现状

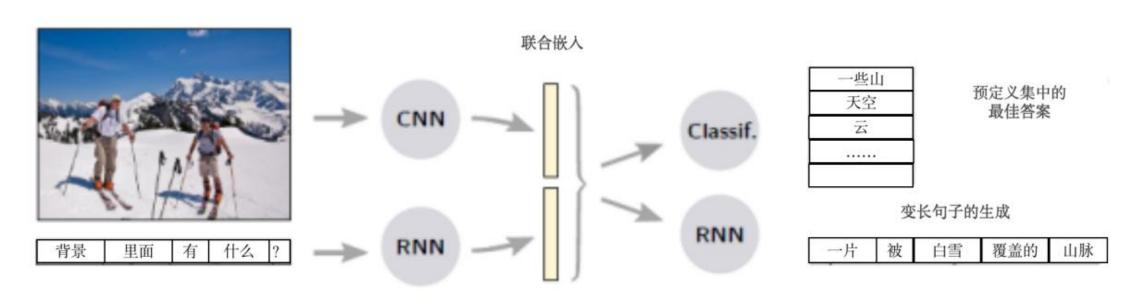


给定一张图片和一个与该图片相关的自然语言 问题, 计算机能产生一个正确的回答

融合CV与NLP的技术 → 计算机需要同时学会理解图像和文字



#### 融合嵌入方法 → 对图像和问题进行联合编码







→ 注意力机制 → 源于机器翻译问题

➡ 让模型动态地调整对输入项各部分的关注度,从而提升模型的"专注力"

设计一种模块化的模型,可根据问题的类型动态组装模块 ● 复合模型 → 来产生答案



使用先验知识库方法



#### 图像-文本多模态任务举例及研究现状

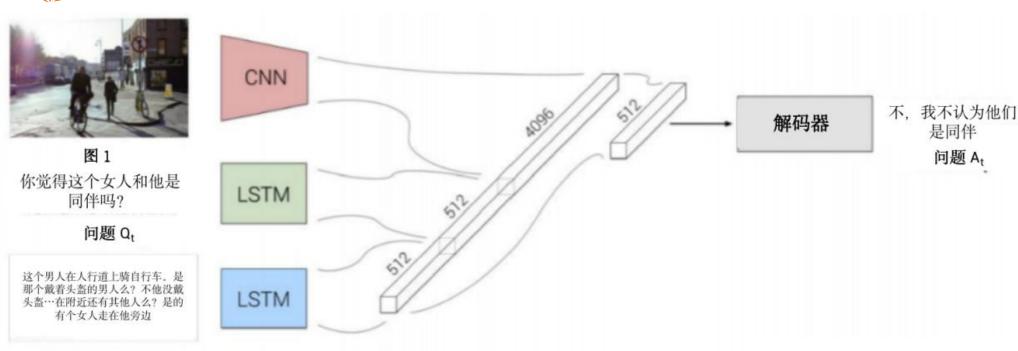


#### ₩ 视觉对话 → AI代理与人类以自然的会话语言对视觉内容进行有意义的对话

#### 具有访问和理解的 需要一个可以组合多个信息源的编码器 多轮对话历史



前t轮信息 (拼接)



# 本章总结

CHAPTER SUMMARY

本章首先介绍了超大规模预训练网络的基本原理及典型网络代表,之后从模型剪枝、模型量化及模型蒸馏三个方面讲解了主流模型压缩方法,然后从理解与生成两个方面,汇总介绍了其它热门研究点并介绍了问答系统与机器阅读理解等任务,最后在多模态学习概念的基础之上,介绍了图像-文本多模态任务的任务定义及网络架构。