

机器翻译前沿综述

冯洋

中科院计算所

2019.10.20



中国科学院计算技术研究所
INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES

ACL 2017

- Translation
- Machine
- Neural
- **Domain**
- Data
- **Adaptation**
- Decoding
- Learning
- **Monolingual**
- Source
- Attention
- Models
- Syntax
- Context
- System

ACL 2018

- Translation
- Machine
- Neural
- Learning
- Attention
- **Unsupervised**
- Multi
- Context
- Model
- Evaluation
- **Document**
- **Non-autoregressive**
- Decoder
- Search

ACL 2019

- Neural
- Learning
- Language
- Translation
- Machine
- Word
- Sentence
- **Embedding**
- Representation
- Multi
- **Unsupervised**
- Semantic
- Attention
- **Multilingual**
- **Adversarial**
- Context

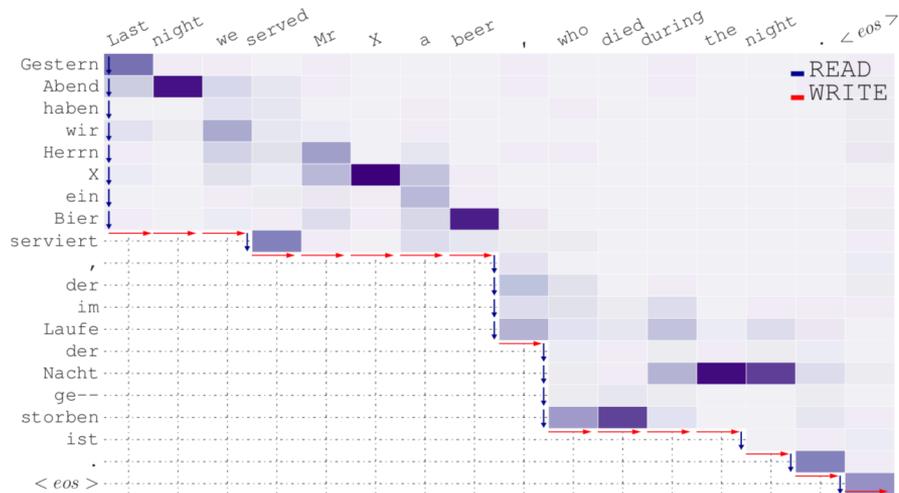
目录

CONTENTS

- 一 同声传译
- 二 多模态机器翻译
- 三 非自回归模型
- 四 篇章翻译
- 五 领域自适应
- 六 多语言翻译
- 七 模型训练

同声传译

- 场景：
- 在读取源语言输入的同时完成翻译的输出，翻译过程仅有几个单词的延迟。



同声传译

- 要解决的问题：
 - 读取源语言输入的过程中，判断是否在当前位置进行输出。
 - SOV语言到SVO语言的翻译，需要读取源语言动词才能进行准确翻译。
 - 如何在训练时模拟测试时的增量环境。

同声传译

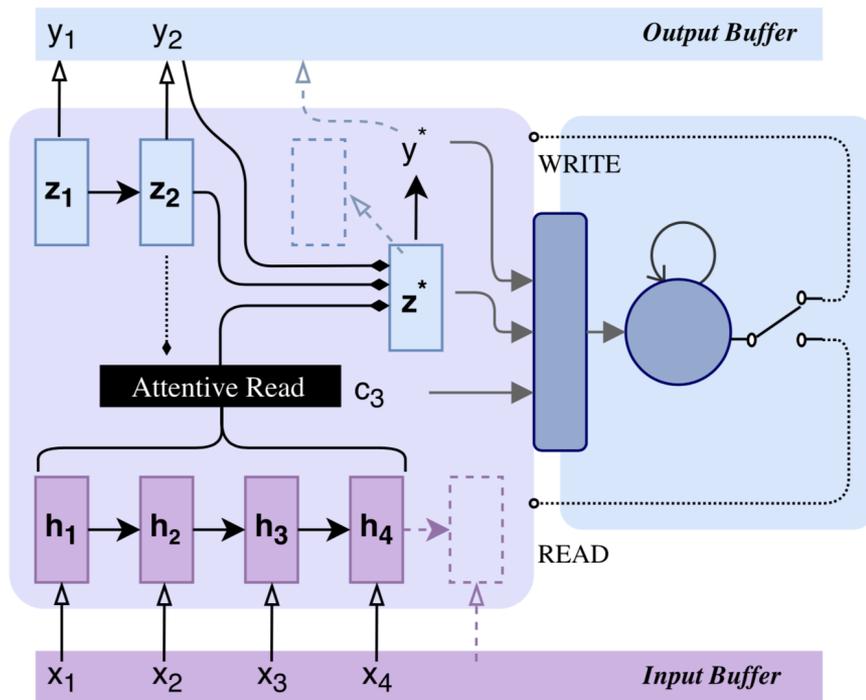
- 解决方案：
 - 更精确的读写策略。

同声传译

- 解决方案：
 - 更精确的读写策略。发展分为三个阶段：
 - 基于解码策略和强化学习

同声传译

- 解决方案
- 更精确
- 基于



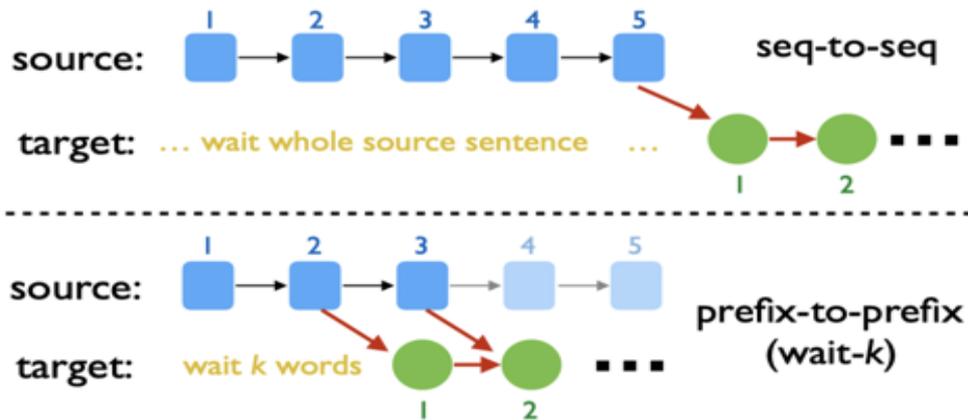
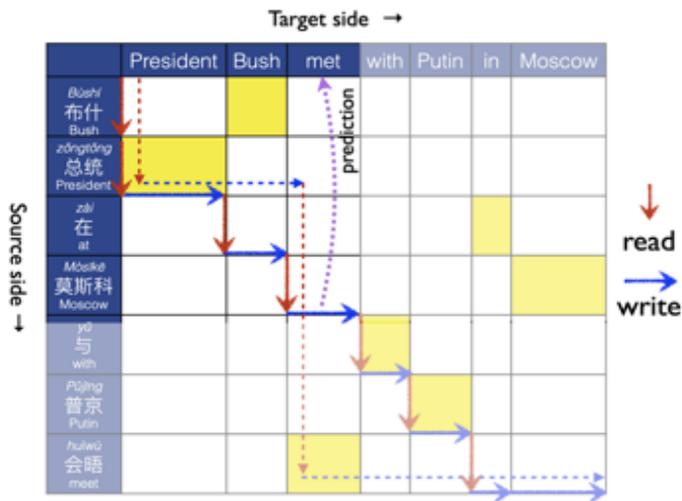
阶段：

每个输入位置都对应读、写两个动作，为每个动作都定义了reward，通过强化学习来使得NMT系统选择reward最大的动作序列。

同声传译

- 解决方案：
 - 更精确的读写策略。发展分为三个阶段：
 - 基于解码策略和强化学习
 - 固定延迟的读写策略，如 wait-k

同声传译



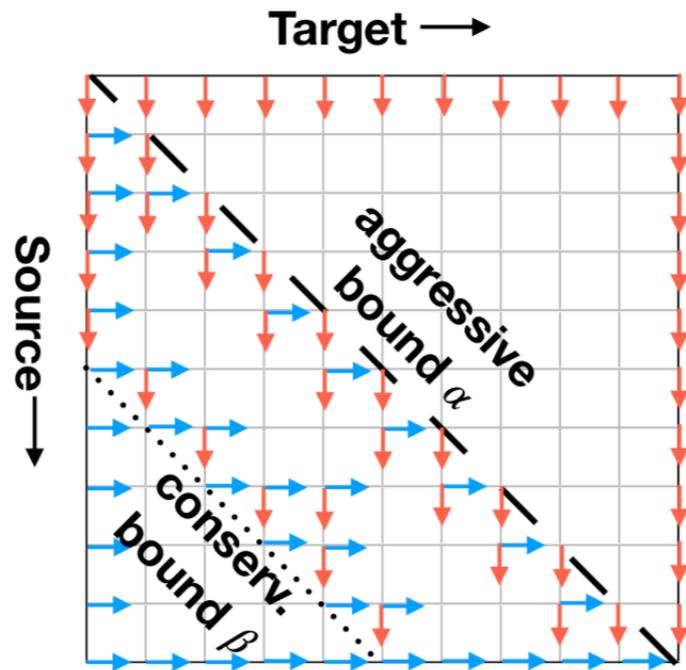
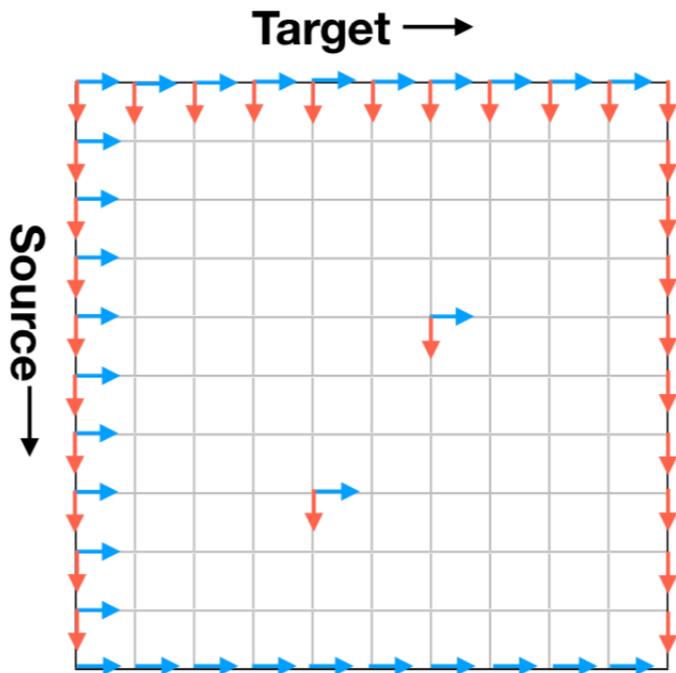
等待输入k个源端单词，然后再进行翻译，即输出始终落后输入k个词。训练时，通过mask源端单词来模拟增量翻译。

同声传译

- 解决方案：
 - 更精确的读写策略。发展分为三个阶段：
 - 策略解码
 - 固定延迟的读写策略，如 wait-k
 - 自适应的读写策略

同声传译

- 解



同声传译

- 解决方案：
 - 更精确的读写策略。
 - 增加模型预测能力。
 - 基于语言模型

多模态翻译

- 场景：
 - 输入为文本、语音、图像、视频中等，经过机器翻译，输出为文字。

多模态机器翻译

- 发展历程：
 - 语音翻译

多模态机器翻译

- 发展历程：
 - 语音翻译
 - 管道模型

多模态机器翻译

- 发展历程：

- 语音翻译

- 管道模型

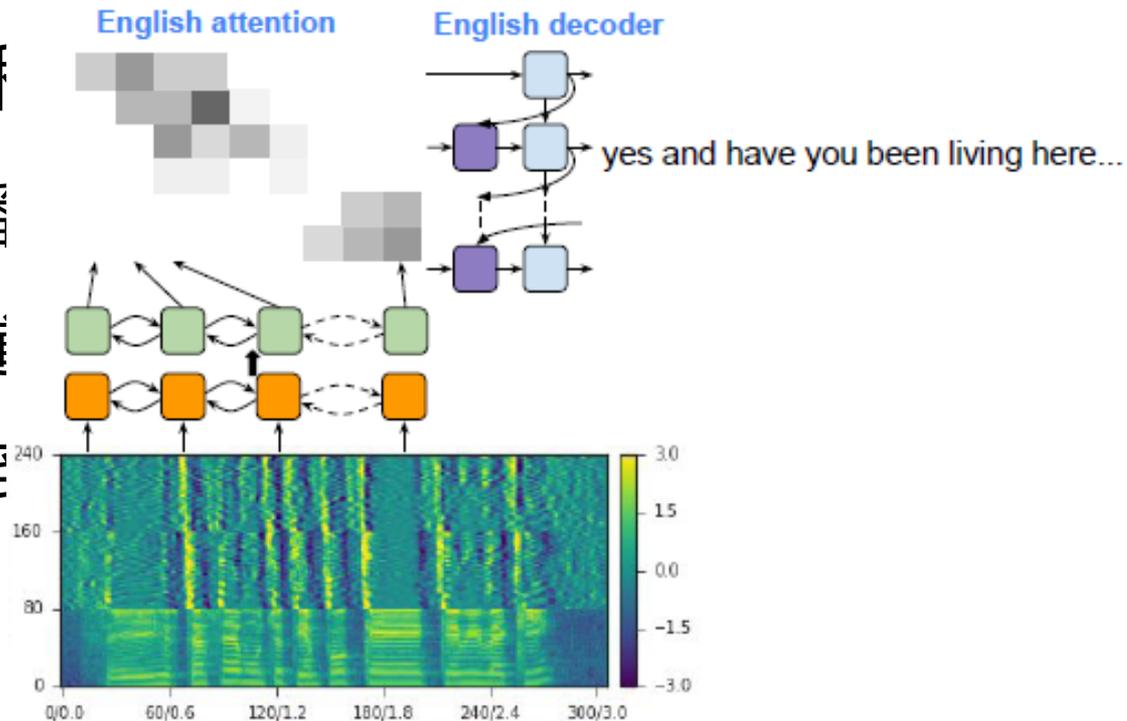


多模态机器翻译

- 发展历程：
 - 语音翻译
 - 管道模型
 - 端到端模型

多模态机器翻译

- 发展历程
- 语音翻译
- 管道
- 端到端



多模态机器翻译

- 发展历程：
 - 语音翻译
 - 管道模型
 - 端到端模型
 - 二者融合，双注意力模型

多模态机器翻译

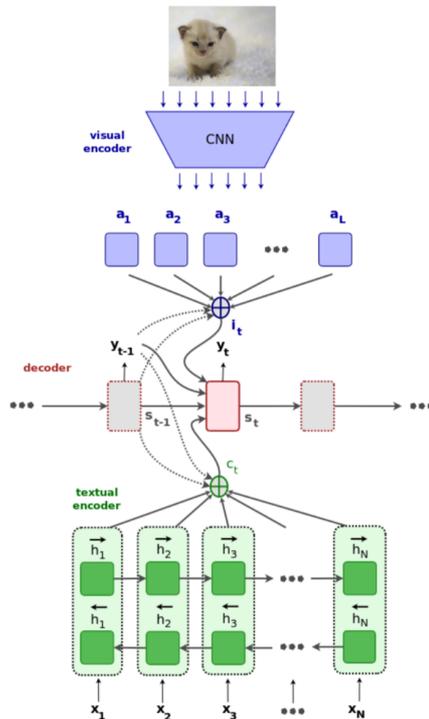
- 发展历程：
 - 语音翻译
 - 多模态翻译

多模态机器翻译

- 发展历程：
 - 语音翻译
 - 多模态翻译
 - 直接使用图像信息

多模态机器翻译

- 发展历程：
 - 语音翻译
 - 多模态翻译
 - 直接使用图像



多模态机器翻译

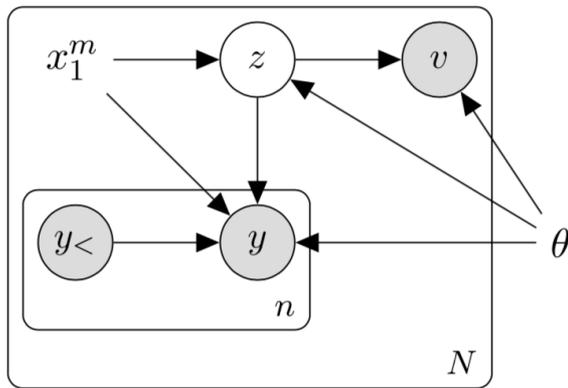
- 发展历程：
 - 语音翻译
 - 多模态翻译
 - 直接使用图像信息
 - 使用Caption信息

多模态机器翻译

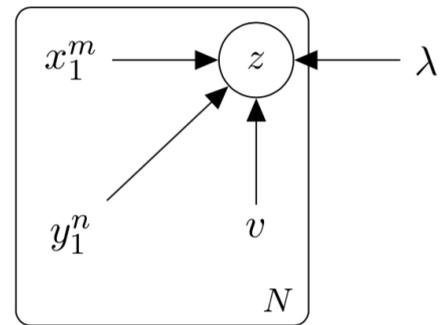
- 发展历程：
 - 语音翻译
 - 多模态翻译
 - 直接使用图像信息
 - 使用Caption信息
 - 使用图像语义

多模态机器翻译

4. 模型



(a) VMMT_C: given the source text x_1^m , we model the joint likelihood of the translation y_1^n , the image (features) v , and a stochastic embedding z sampled from a conditional latent Gaussian model. Note that the stochastic embedding is the sole responsible for assigning a probability to the observation v , and it helps assign a probability to the translation.



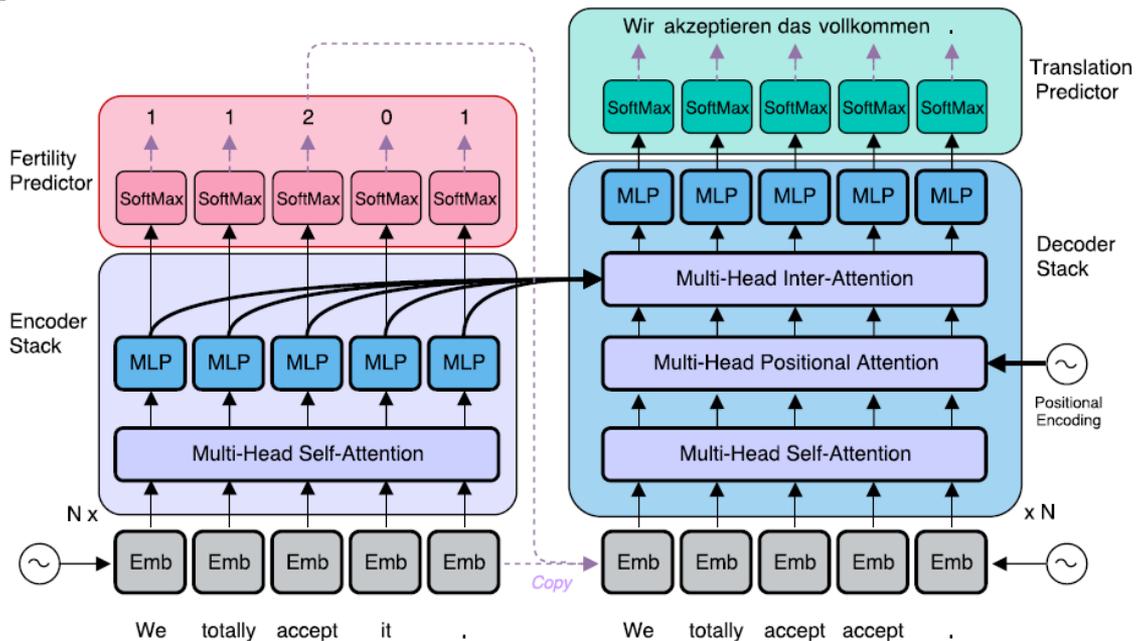
(b) Inference model for VMMT_C: to approximate the true posterior we have access to both modalities (text x_1^m , y_1^n and image v).

非自回归模型

- 场景：
 - 所有目标词同时生成，无需将前一个生成词作为输入

非自回归模型

通用模型：

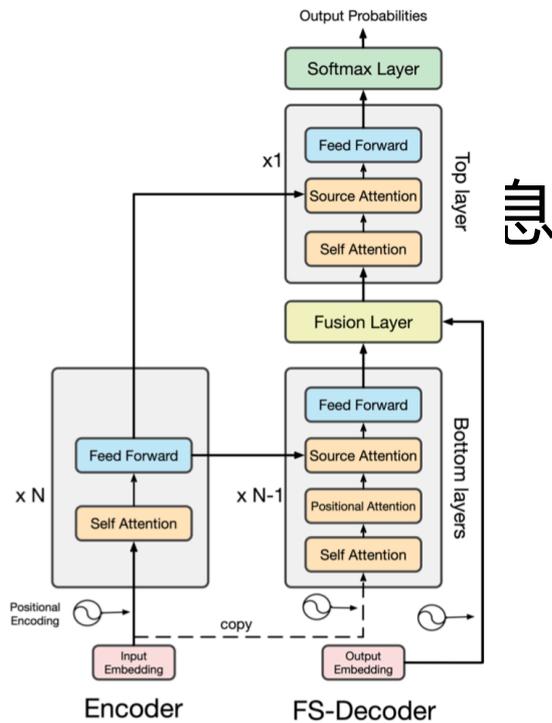


非自回归模型

- 发展趋势：添加序列信息
 - 方式一：在模型中添加序列信息

非自回归模型

- 发展趋势：添力
- 方式一：在模



非自回归模型

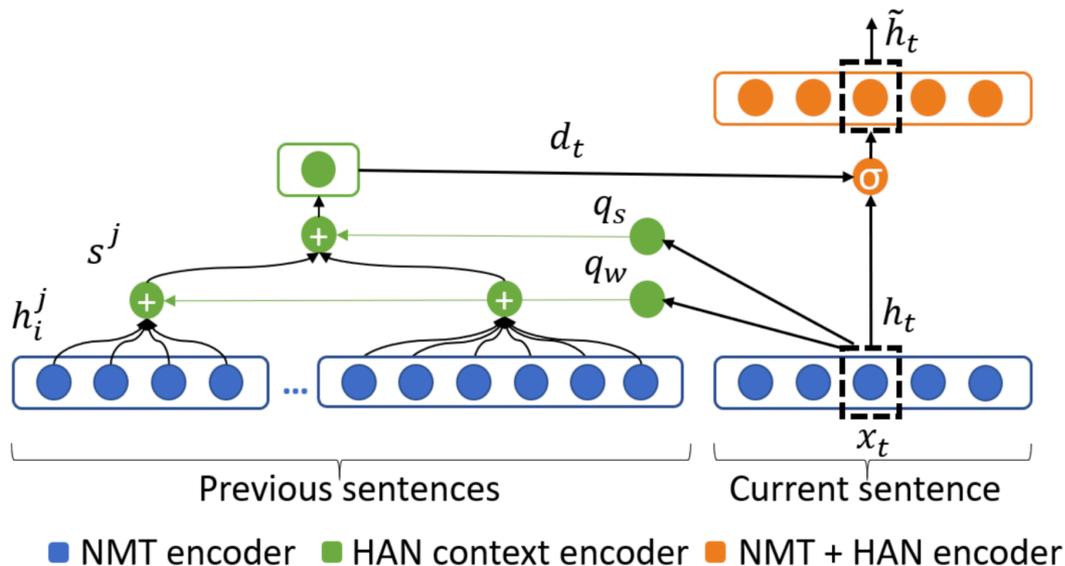
- 发展趋势：添加序列信息
 - 方式一：在模型中添加序列信息
 - 方式二：序列级训练，如强化学习

篇章翻译

- 场景：
 - 对一个段落进行翻译，需考虑上下文的一致性以及指代问题

篇章翻译

• 通用模型：



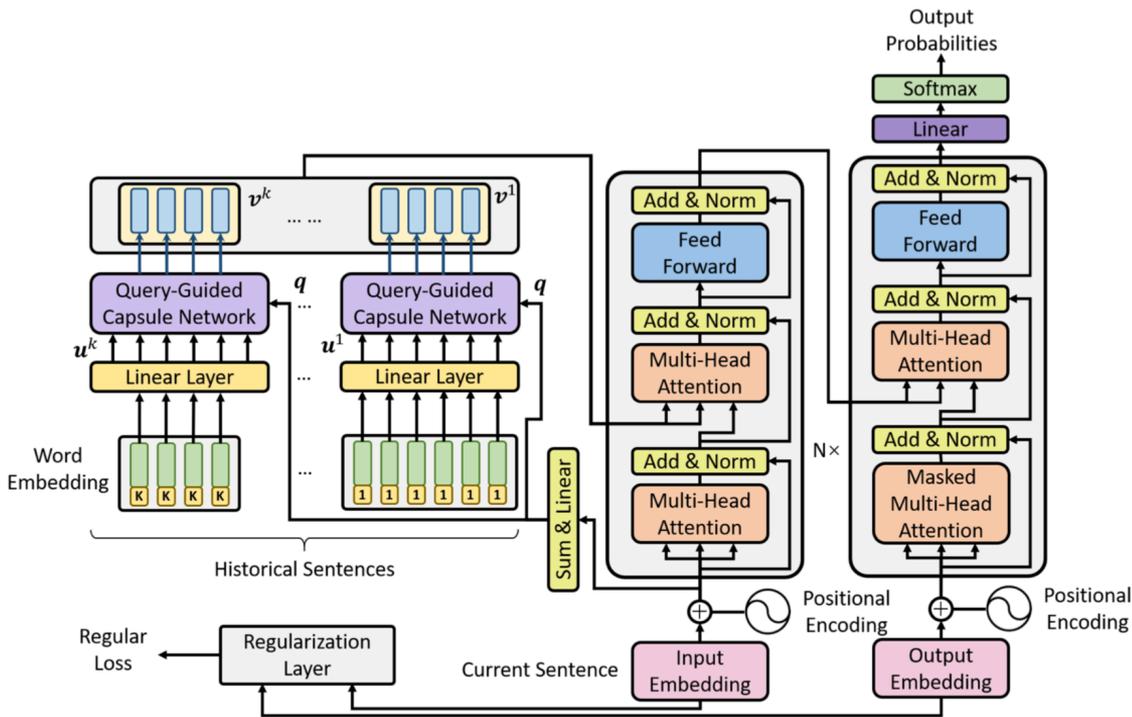
篇章翻译

- 发展趋势：
 - 改进注意力机制

篇章翻译

• 发展趋势

• 改进注



篇章翻译

- 发展趋势：
 - 改进注意力机制
 - 学习更好的表示 (参照多轮对话)

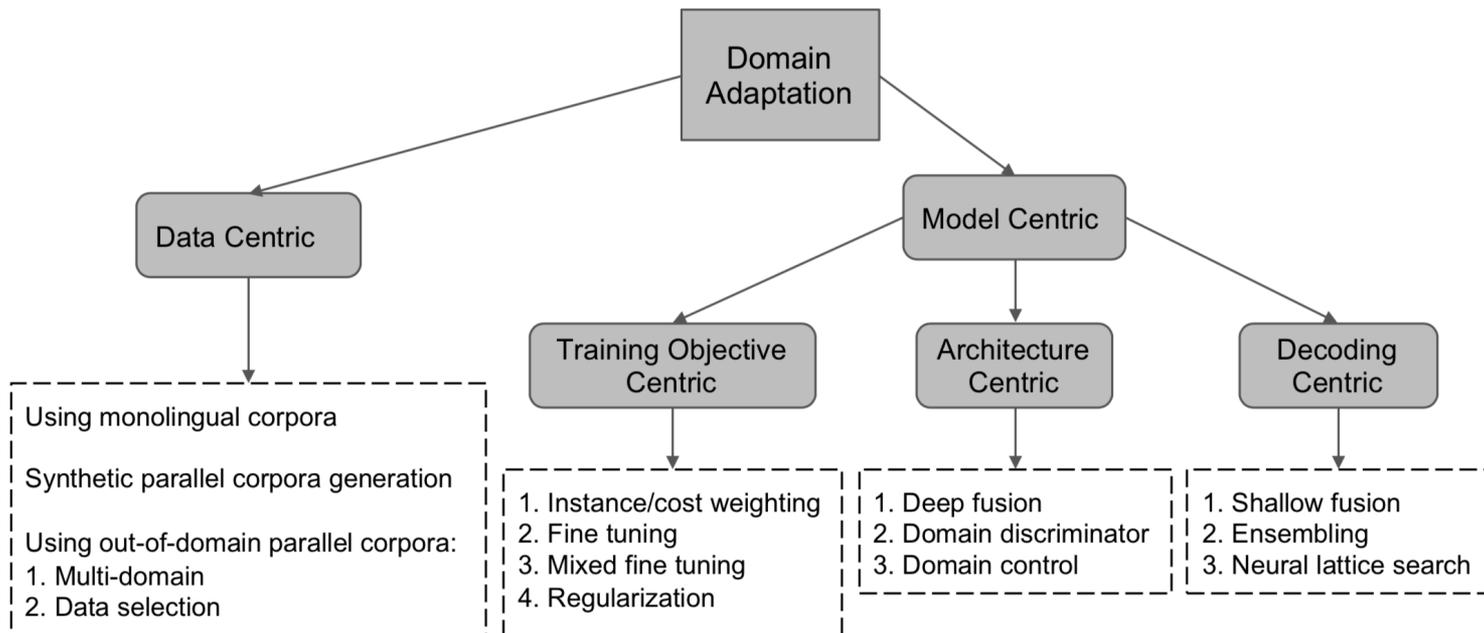
篇章翻译

- 发展趋势：
 - 改进注意力机制
 - 学习更好的表示
 - 使用更长的历史

领域自适应

- 场景：
 - 大量的out-of-domain或者通用领域数据，少量的in-domain数据，需要对in-domain的句子进行翻译。

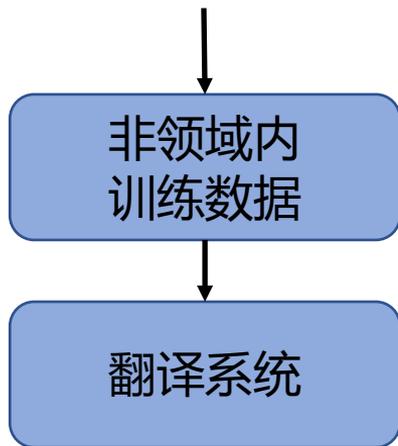
领域自适应



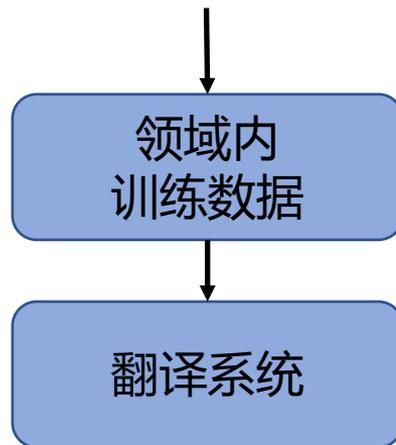
领域自适应

- 通用模型：Fine tune

第一阶段：



第二阶段：

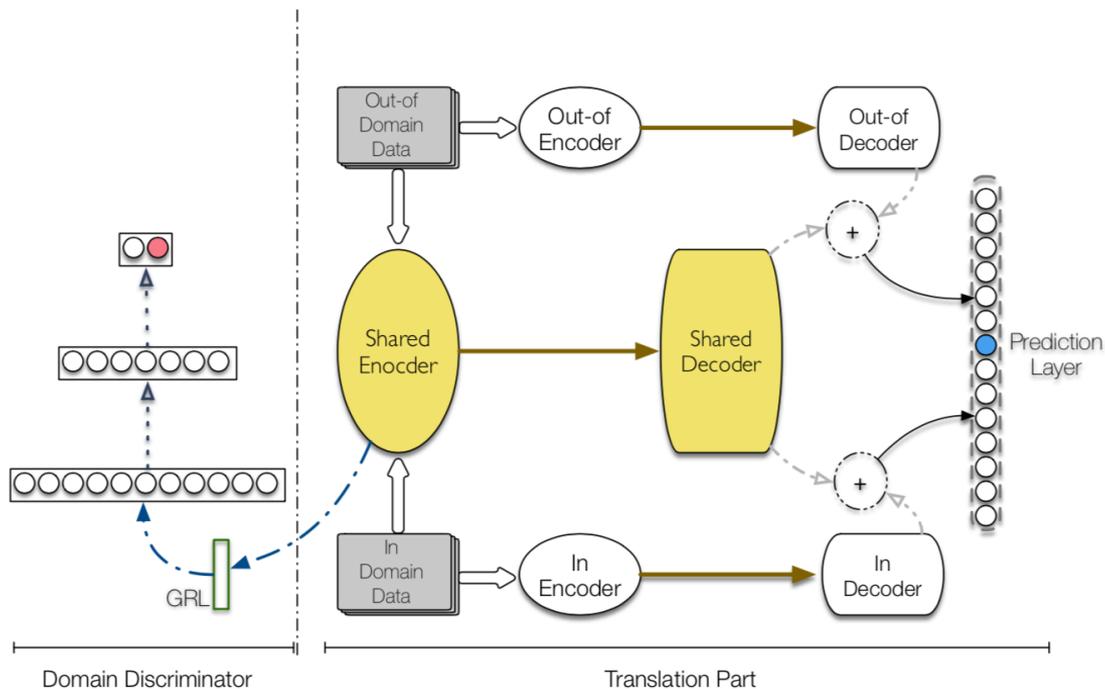


领域自适应

- 发展趋势：
 - 采用对抗学习抽取特征

领域自适应

- 发展趋势
- 采用双



领域自适应

- 发展趋势：
 - 采用对抗学习抽取特征
 - 无监督领域自适应

领域自适应

- 发展趋势：
 - 采用对抗学习抽取特征
 - 无监督领域自适应
 - 课程学习进行数据选择

多语言翻译

- 场景：
 - 一对多
 - 多对一
 - 多对多

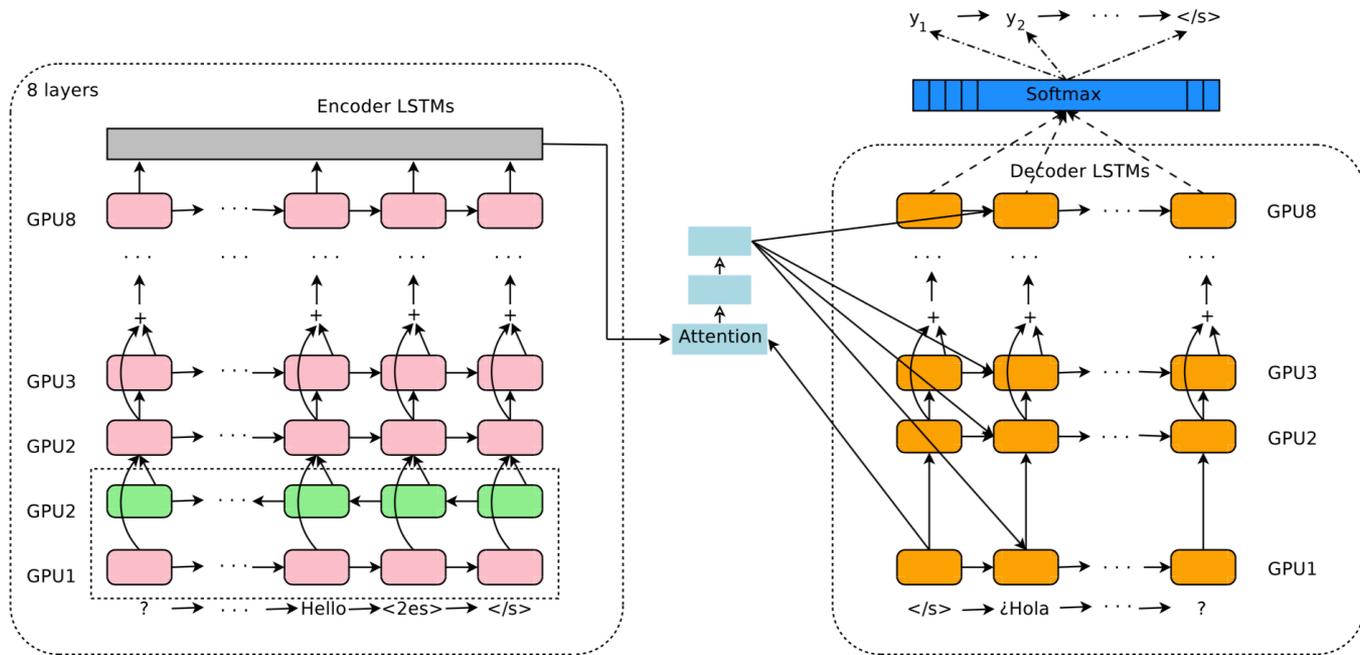
多语言翻译

- 发展趋势：
 - 共享结构

多语言翻译

• 发展

• 共



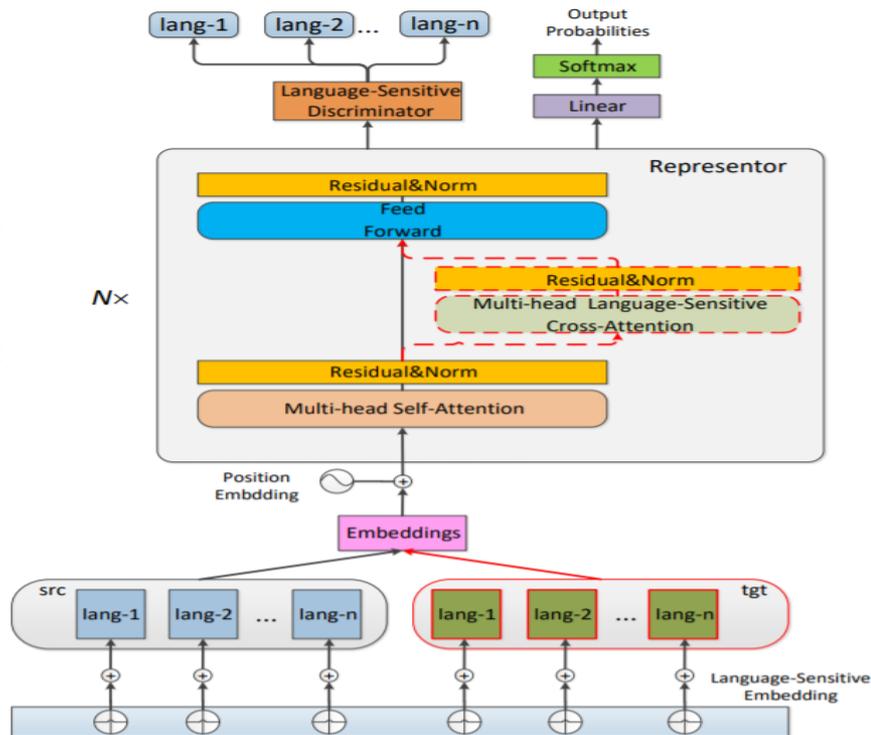
Melvin Johnson, Mike Schuster, Quoc V Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, et al. 2016. Google's multilingual neural machine translation system: Enabling zero-shot translation.

多语言翻译

- 发展趋势：
 - 共享编码器-解码器
 - 区分语言特征

多语言翻译

- 发展趋势：
 - 共享编码器
 - 区分语言特



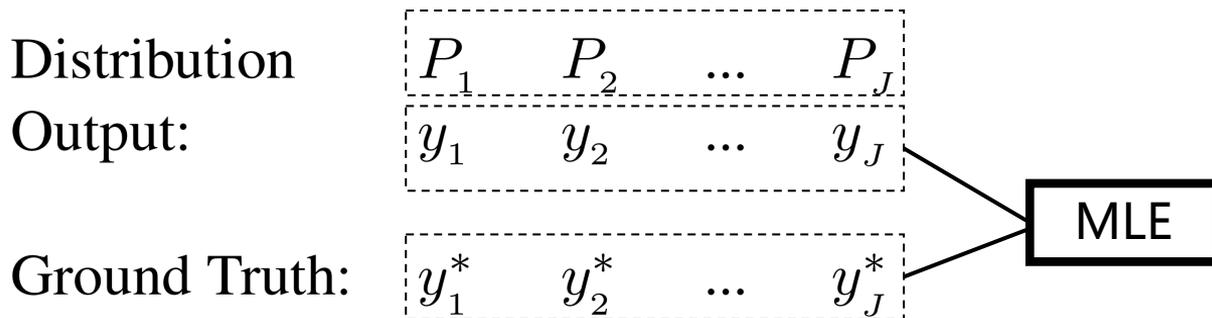
Yining Wang, Long Zhou, Jiajun Zhang, Feifei Zhai, Jingfang Xu and Chengqing Zong. A Compact and Language-Sensitive Multilingual Translation Method. *In Proceedings of ACL 2019.*

多语言翻译

- 发展趋势：
 - 共享编码器-解码器
 - 区分语言特征
 - 改善低资源和零资源翻译

模型训练

- 通用框架：Teacher Forcing



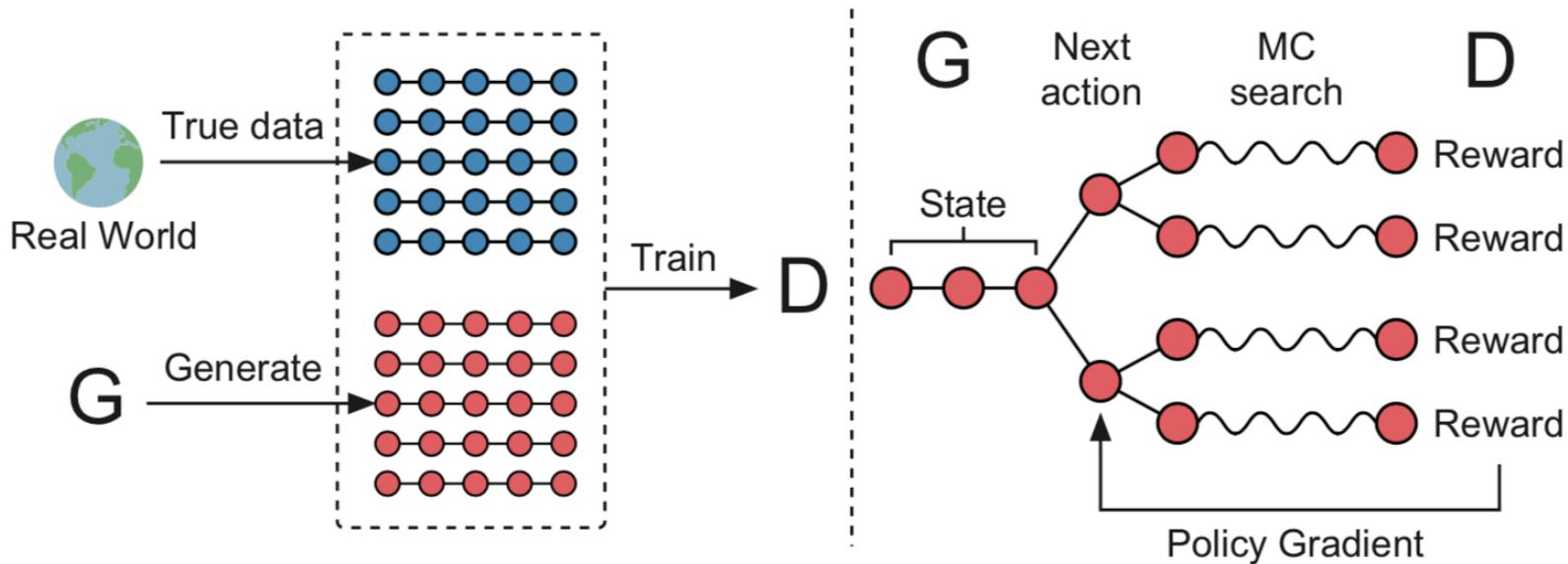
Loss:
$$\mathcal{L}_{\text{MLE}} = - \sum_{j=1}^J \log P_j(y_j^*)$$

模型训练

- 发展趋势：
 - 基于强化学习的序列级训练

模型训练

- 发展趋势：



模型训练

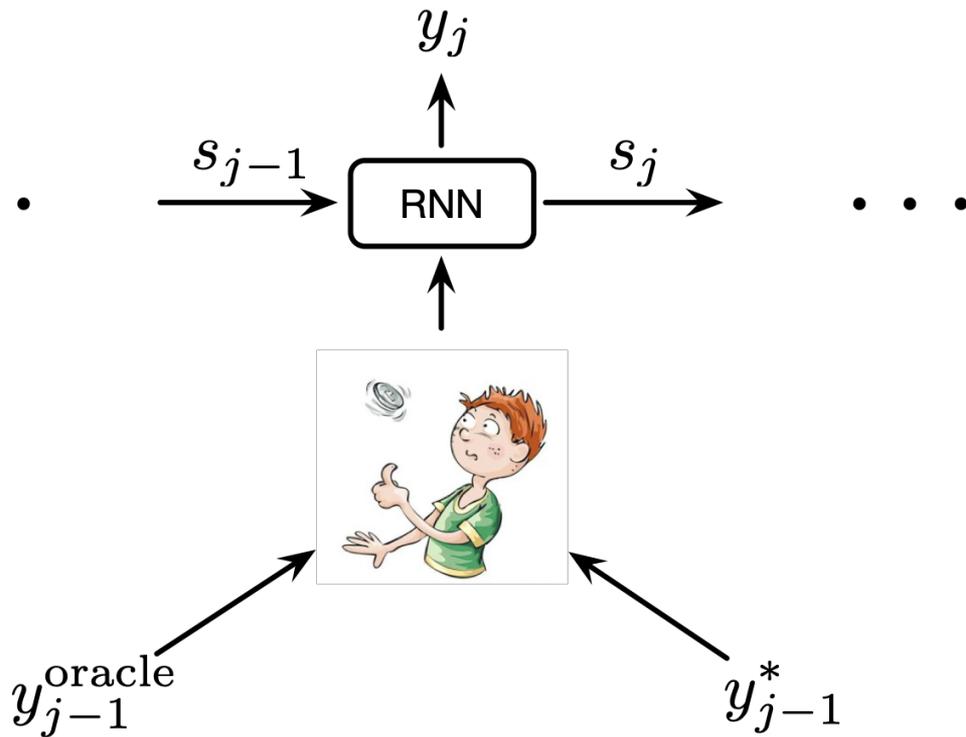
- 发展趋势：
 - 基于强化学习的序列级训练
 - 训练和测试一致性

模型训练

- 发展趋势

- 基于强 ..

- 训练和



模型训练

- 发展趋势：
 - 基于强化学习的序列级训练
 - 训练和测试一致性
 - 可导的序列级训练

低资源翻译

- 通用模型：数据增强，如Back Translation
- 趋势：仍以数据增强为主。
- 最新论文：
 - ACL 2019：
 - [Revisiting](#) Low-Resource Neural Machine Translation: A Case Study
 - Generalized [Data Augmentation](#) for Low-Resource Translation (短文)
 - EMNLP 2019:
 - Handling [Syntactic Divergence](#) in Low-resource Machine Translation
 - Exploiting Multilingualism through Multistage [Fine-Tuning](#) for Low-Resource Neural Machine Translation

Pre-train

- 在机器翻译上没有显著提升
 - 机器翻译语料规模已足够？

目录

CONTENTS

一 同声传译

二 多模态机器翻译

三 非自回归模型

四 篇章翻译

五 领域自适应

六 多语言翻译

七 模型训练

Thanks for your attention!