华为机器翻译架构和模型加速

魏代猛 华为2012 / 机器翻译实验室







精彩继续! 更多一线大厂前沿技术案例

上海站



时间: 2023年4月21-22日

地点:上海·明捷万丽酒店

扫码查看大会详情>>



广州站

全球软件开发大会

时间: 2023年5月26-27日

地点:广州·粤海喜来登酒店

扫码查看大会详情>>



个人简介



- 华为高级技术专家,机器翻译算法负责人,产品落地华为云、HMS、华为手机等
- 北京大学硕士, 研究方向: 机器翻译、同传翻译、语义理解等
- 带领团队参加 WMT20/21/22 news、biomedical、efficiency等赛道多项 第一, IWSLT 22 多项第一,WAT20比赛多项第一
- 在AAAI, ACL, EMNLP, ICASSP等发表论文30+







大纲

- 机器翻译简介
- 模型推理问题
- 端测推理加速
- 华为机器翻译
- 总结

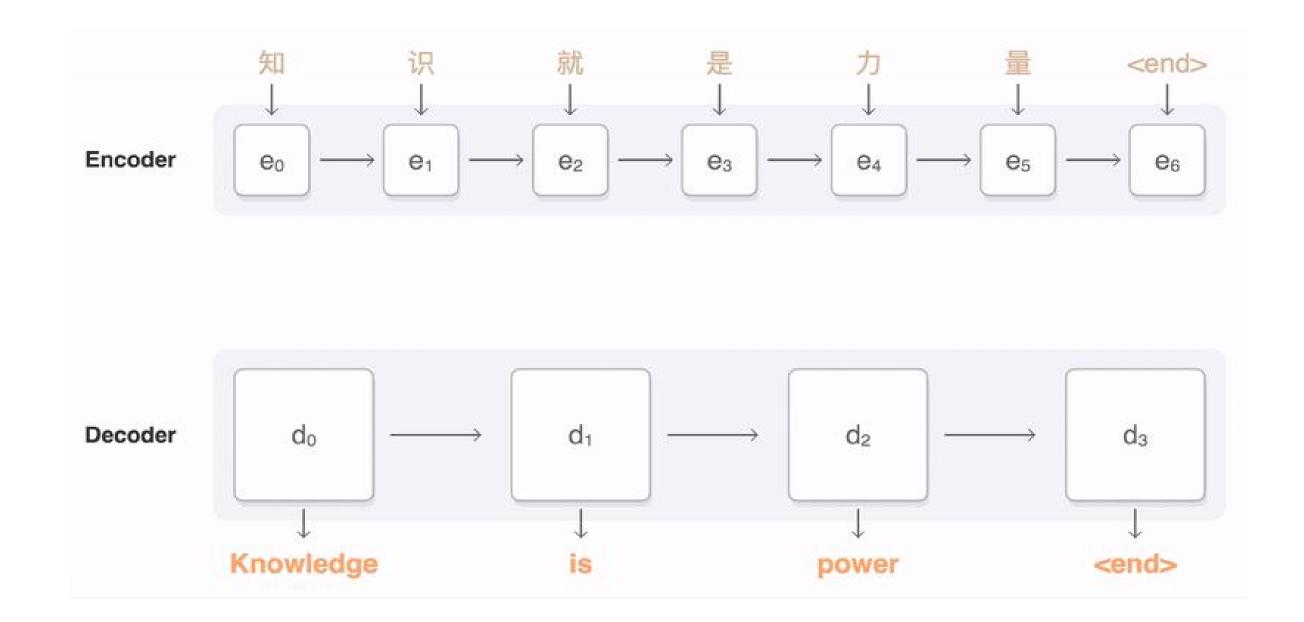






机器翻译简介

- 1、主流的机器翻译模型包含Encoder和Decoder两部分,Encoder将原文整个序列编码成一个多维向量,Decoder将原文序列的向 量解码成译文。
- 2、Attention模型记录原文和译文的词对齐关系,指导机器翻译在解码译文某个词时,应该更关注与原文的哪一个部分,以提升长 句翻译质量。

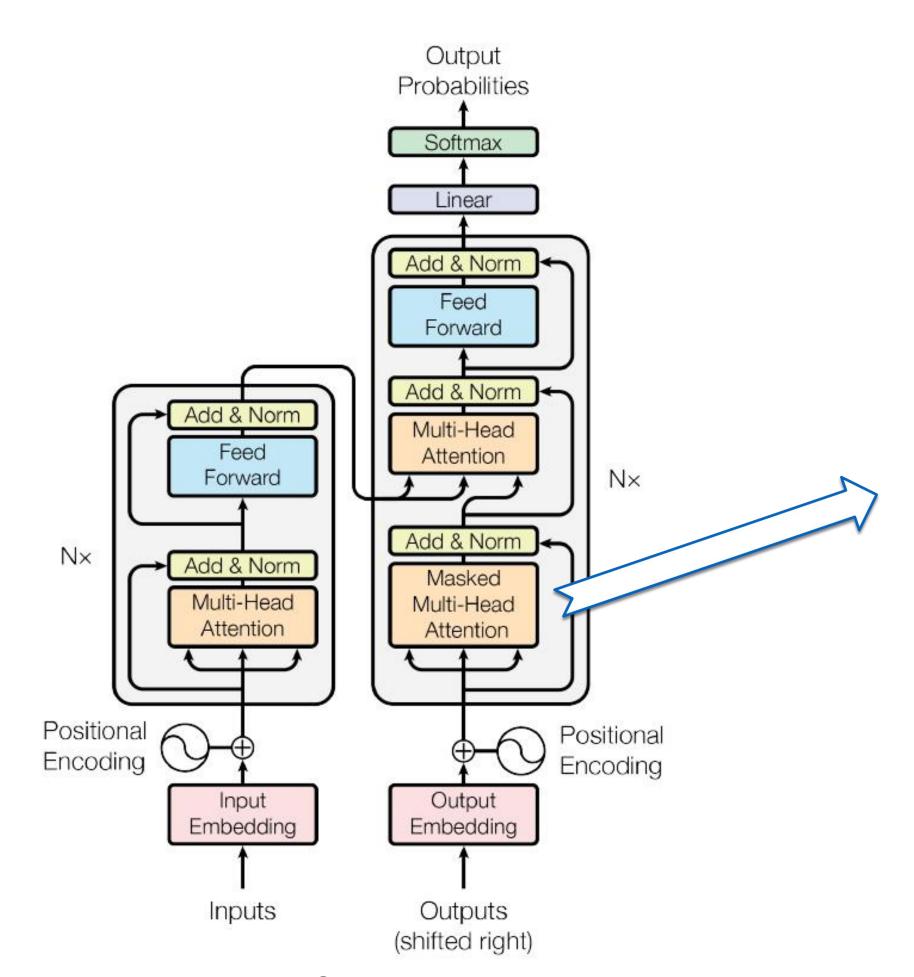








机器翻译简介



Encoder: 6

Decoder: 6

Hidden size: 1024

参数量: 2亿

大小: 800M

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser and Illia Polosukhin. "Attention is All you Need" Neural Information Processing Systems (2017).







Transformer模型在GPU, CPU, ARM运行典型值

	GPU/T4	CPU/Intel	ARM
耗时	45 ms/token	150 ms/token	-

模型太大, 计算量太大 端侧最具挑战

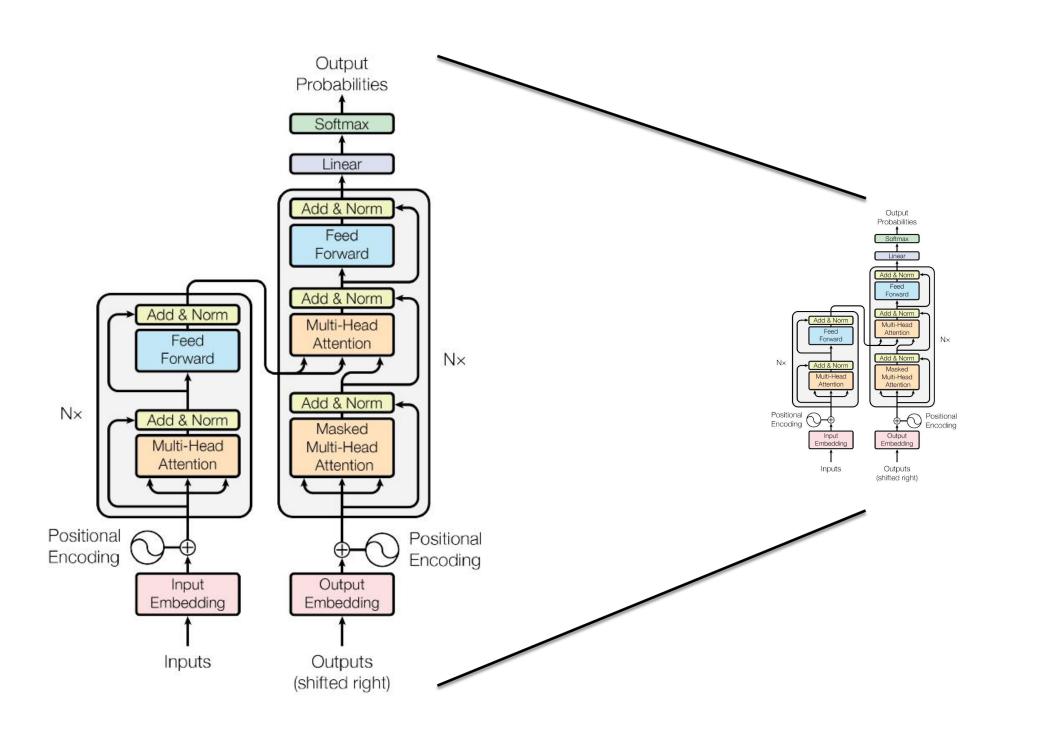






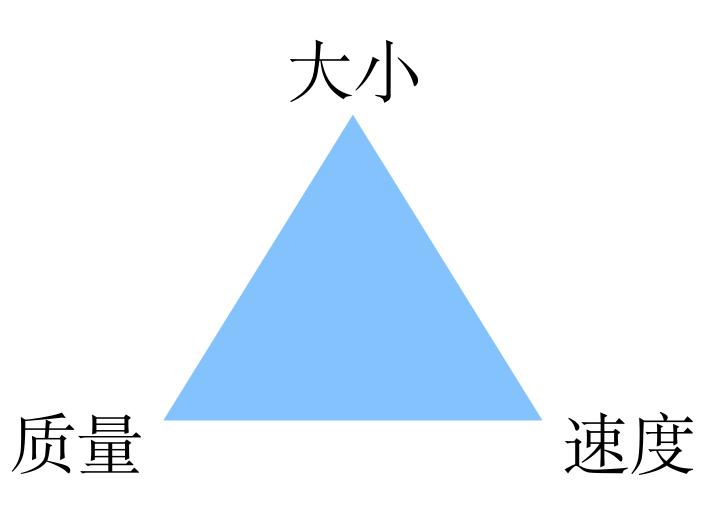
模型推理问题

让模型变小



存储 计算量

质量↓









大纲

- 机器翻译简介
- 模型推理问题
- 端测推理加速
- 华为机器翻译
- 总结

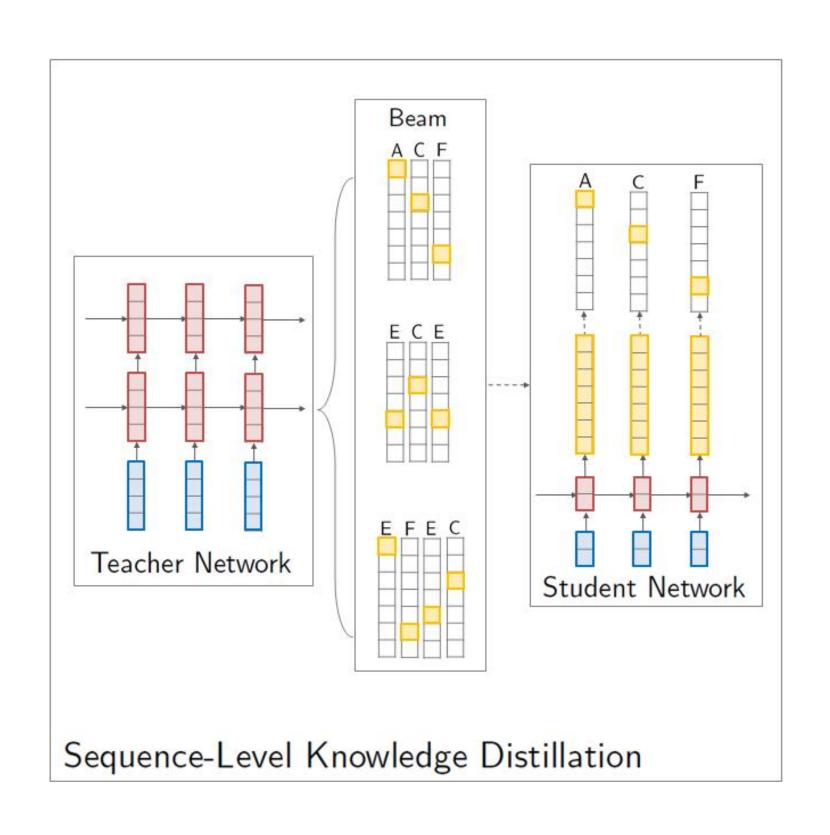






模型存储小 小模型 → 质量变差 计算量小 小模型 → 高质量?

知识蒸馏



Geoffrey E. Hinton, Oriol Vinyals and Jeffrey Dean. "<u>Distilling the Knowledge in a Neural Network</u>" arXiv: Machine Learning (2015): n. pag. Yoon Kim and Alexander M. Rush. "<u>Sequence-Level Knowledge Distillation</u>" Empirical Methods in Natural Language Processing (2016). Markus Freitag, Yaser Al-Onaizan and Baskaran Sankaran. "<u>Ensemble Distillation for Neural Machine Translation</u>" arXiv: Computation and Language (2017): n. pag.







知识蒸馏

Model	Emb.	FFN	Head	Depth	Params(M)	Size(MB)	wmt19	wmt20
Teacher*4	1024	4096	16	25/6	514	2000	46.71	39.70
Base.12	512	2048	8	12/1	53	210	44.65	38.02

96%

小模型&高质量

Tiny Bert

System	#Params	#FLOPs	Speedup	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Avg
BERT _{BASE} (Teacher)	109M	22.5B	1.0x	83.9/83.4	71.1	90.9	93.4	52.8	85.2	87.5	67.0	79.5
BERTTINY	14.5M	1.2B	9.4x	75.4/74.9	66.5	84.8	87.6	19.5	77.1	83.2	62.6	70.2
$BERT_{SMALL}$	29.2M	3.4B	5.7x	77.6/77.0	68.1	86.4	89.7	27.8	77.0	83.4	61.8	72.1
BERT ₄ -PKD	52.2M	7.6B	3.0x	79.9/79.3	70.2	85.1	89.4	24.8	79.8	82.6	62.3	72.6
DistilBERT ₄	52.2M	7.6B	3.0x	78.9/78.0	68.5	85.2	91.4	32.8	76.1	82.4	54.1	71.9
MobileBERT _{TINY} †	15.1M	3.1B	<u>=</u>	81.5/81.6	68.9	89.5	91.7	46.7	80.1	87.9	65.1	77.0
TinyBERT ₄ (ours)	14.5M	1.2B	9.4x	82.5/81.8	71.3	87.7	92.6	44.1	80.4	86.4	66.6	77.0

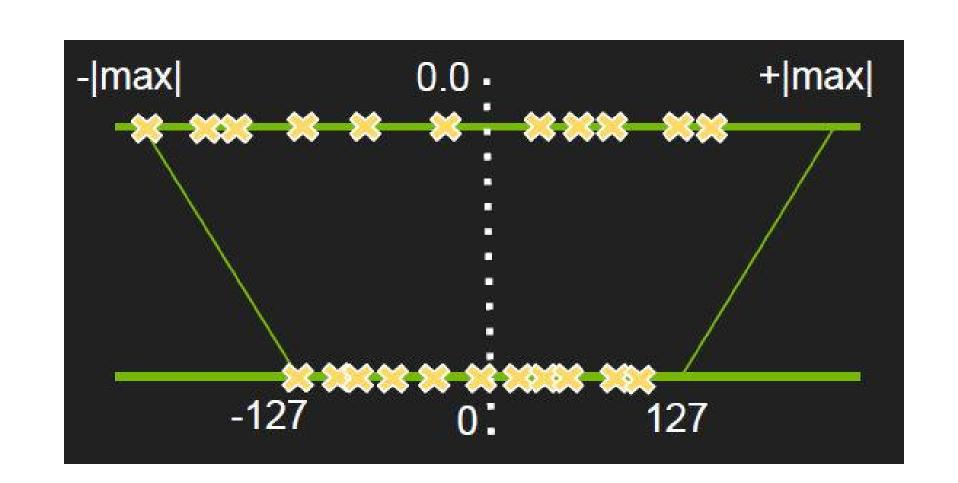
Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. TinyBERT: Distilling BERT for Natural Language Understanding. EMNLP 2020

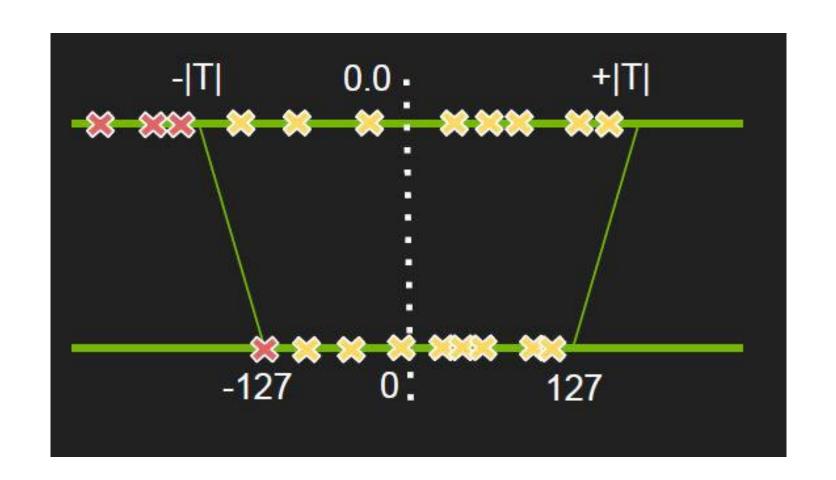






小模型: 更小空间, 更快的推理 -----模型压缩, 低精度推理





增加量化层,FP32->Int8->FP32 , E2E训练

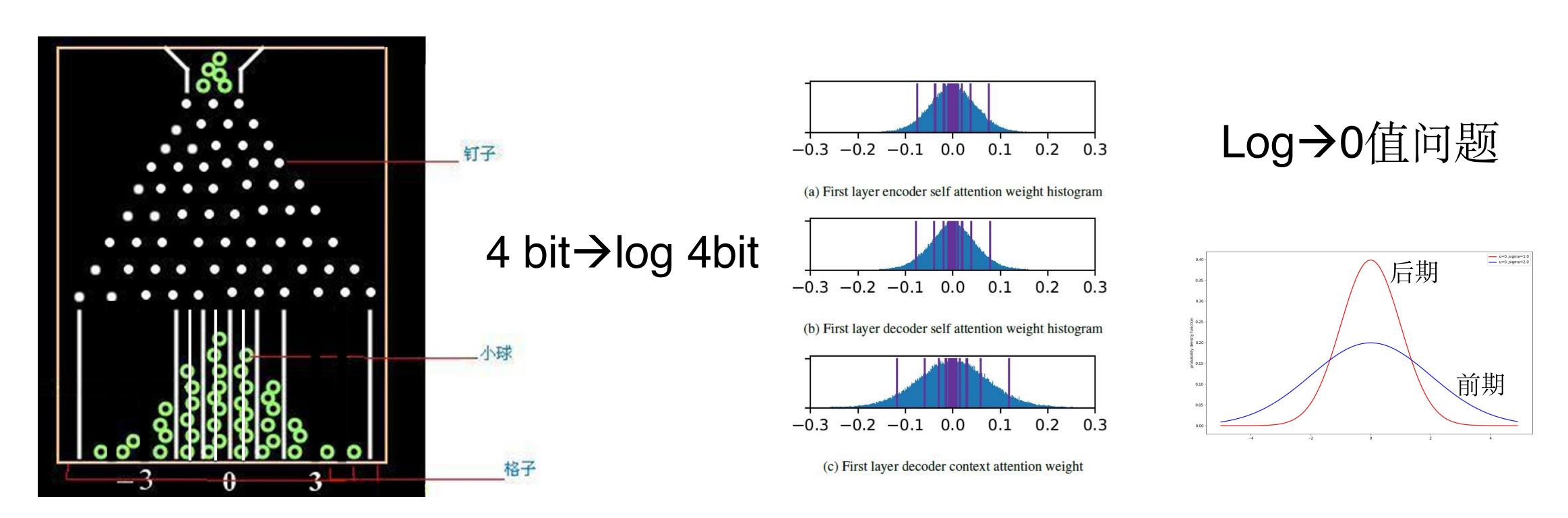
Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam and Dmitry Kalenichenko. "Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference" Computer Vision and Pattern Recognition (2018).







小模型: 更小空间, 更快的推理 -----模型压缩, 低精度推理



Alham Fikri Aji and Kenneth Heafield. 2020. Compressing Neural Machine Translation Models with 4-bit Precision. In *Proceedings of the Fourth Workshop on Neural Generation and Translation*, pages 35–42, Online. Association for Computational Linguistics.







小模型: 更小空间, 更快的推理 -----模型压缩, 低精度推理

模型	BLEU(WMT14)	Parameter Size
32 bit	26.5	260M
8 bit	26.4(-0.1)	66M
4 bit	24.3(-2.2)	35M
Log 4 bit(后期)	25.1(-1.4)	35M
Log 4 bit(前期)	26.2(-0.3)	35M

直接4bit 影响大,log 4 bit前期介入量化训练很关键







小模型: 更小空间, 更快的推理 -----模型压缩, 低精度推理

Int8推理

模型中计算量最大的是矩阵运算(GEMM)

Int8推理: 用整型运算代替浮点型运算提速

处理好量化和反量化是提速的关键

华为Noah高性能推理实验室 https://github.com/huawei-noah/bolt

Operator	FP32	INT8
FC	4.806	1.53
Transpose	0.4158	0.228
Eltwise	0.198	0.13
MatMul	0.1134	0.068
Softmax	0.063	0.040
LayerNorm	0.063	0.036
Reshape	0.027	0.013
Activation	0.0216	0.011
Embedding	0.018	0.010
Slice	0.009	0.003
Multiply	0.0054	0.001
Total	5.94	2.1

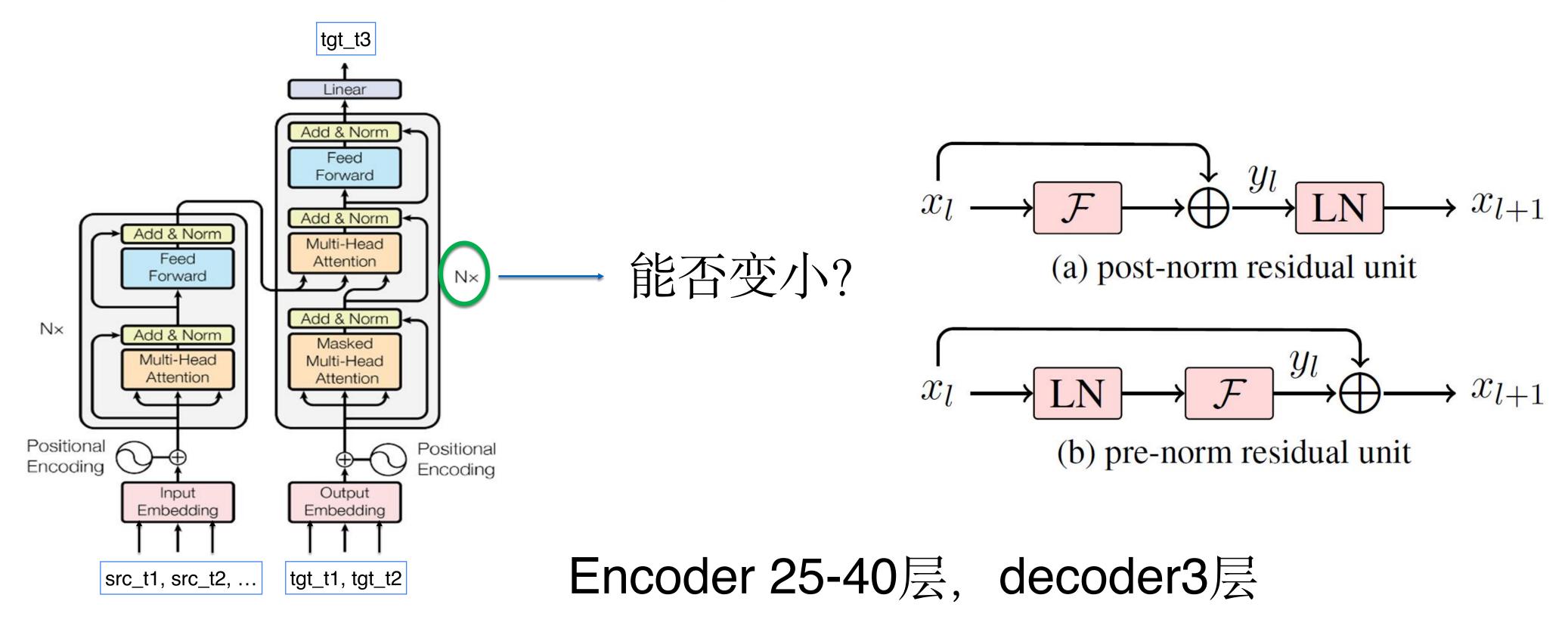
Daya Shanker Khudia, Jianyu Huang, Protonu Basu, Summer Deng, Haixin Liu, Jongsoo Park and Mikhail Smelyanskiy. "FBGEMM: Enabling High-Performance Low-Precision Deep Learning Inference.." arXiv: Learning (2021): n. pag.







小模型: 更强的能力 -----结构优化, 参数共享, 多语言模型



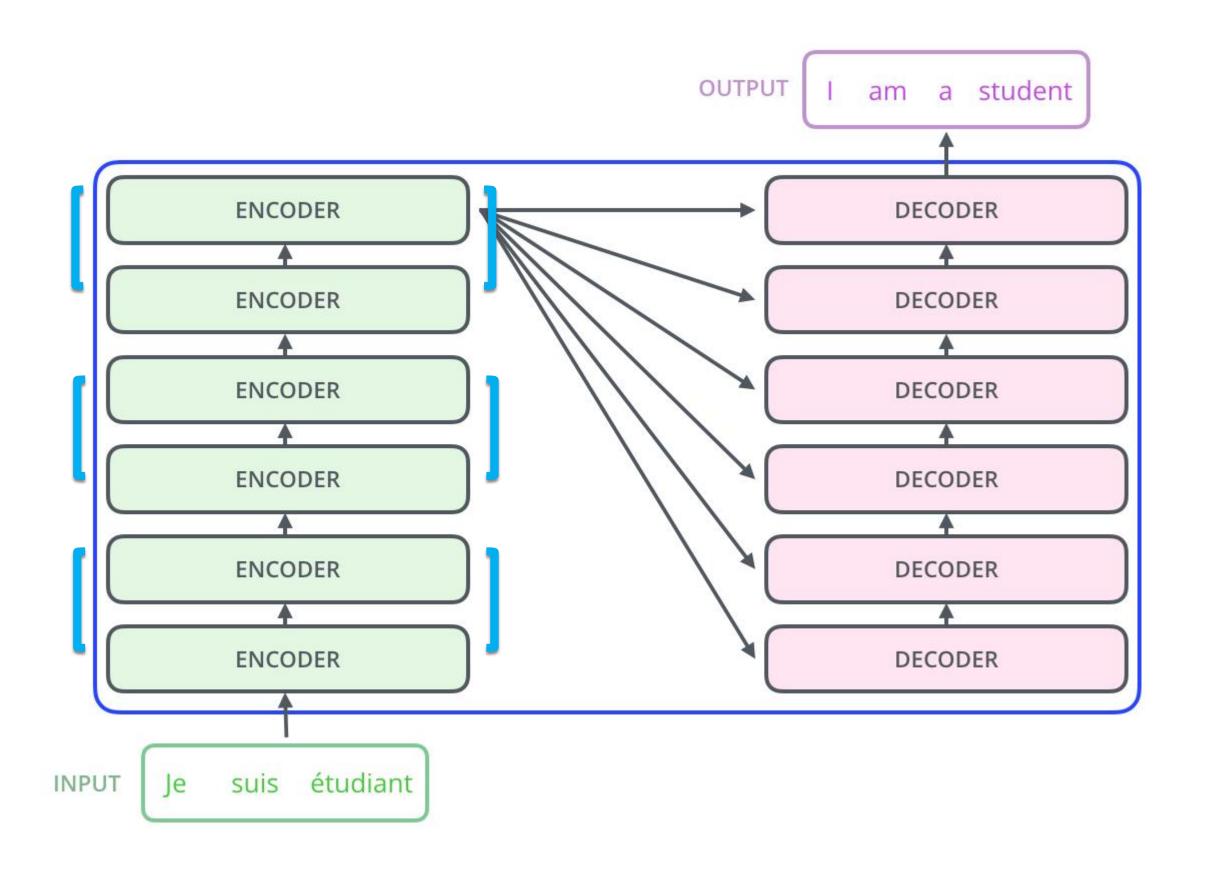
Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F. Wong and Lidia S. Chao. "Learning Deep Transformer Models for Machine Translation.." Meeting of the Association for Computational Linguistics (2019).







小模型: 更强的能力 -----结构优化, 参数共享, 多语言模型



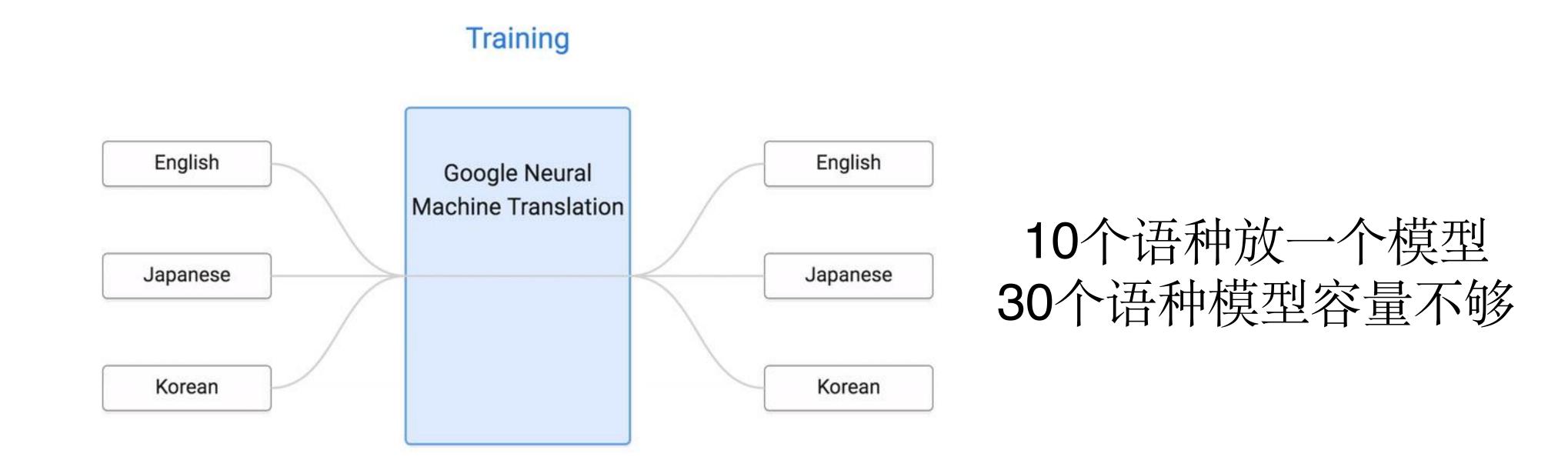
相邻层共享最好







小模型: 更强的能力 -----结构优化,参数共享,多语言模型



Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda B. Viégas, Martin Wattenberg, Greg S. Corrado, Macduff Hughes and Jeffrey Dean. "Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation" Transactions of the Association for Computational Linguistics 5 (2017): 339-351.

Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George Foster, Colin Cherry, Wolfgang Macherey, Zhifeng Chen and Yonghui Wu. "Massively Multilingual Neural Machine Translation in the Wild: Findings and Challenges" arXiv: Computation and Language (2019): n. pag.







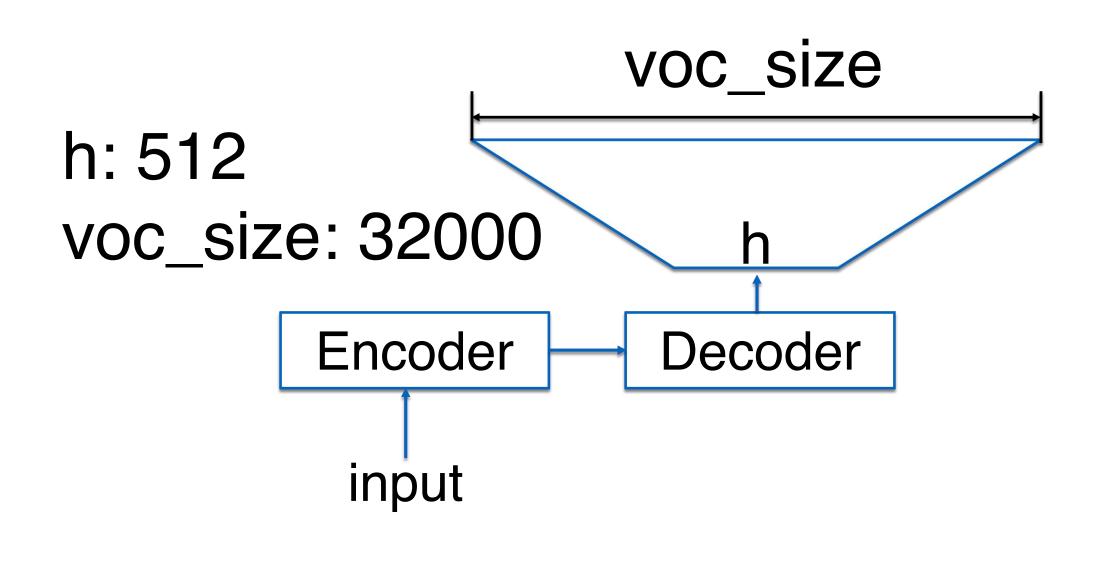
策略	质量	速度	大小
知识蒸馏	$\sqrt{}$	$\sqrt{}$	_
量化推理	_	$\sqrt{}$	$\sqrt{}$
模型结构	$\sqrt{}$	$\sqrt{}$	_
参数共享	$\sqrt{}$	_	\checkmark
多语言	$\sqrt{}$	_	$\sqrt{}$
ShortList		工、定、八、八、上、左、三	
Decoder结构	マ壮-	于减少计算量	



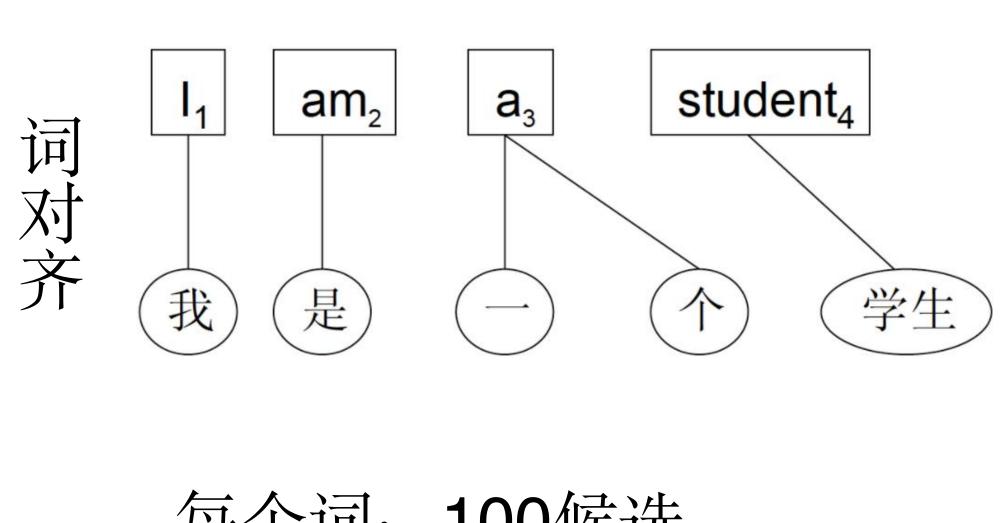




ShortList优化



10 x FFN



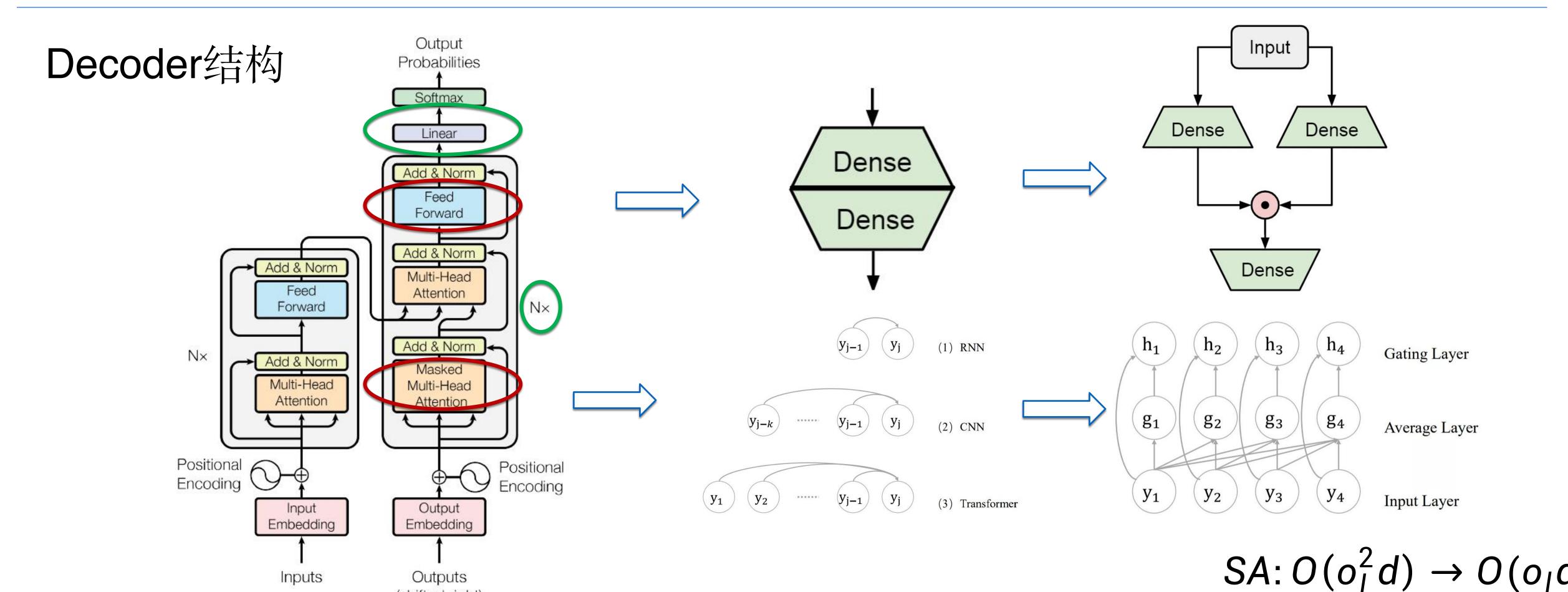
每个词: 100候选 16词→300候选(去重) 512x32000 → 512x300

词对齐可以用Fastalign,每个词75候选









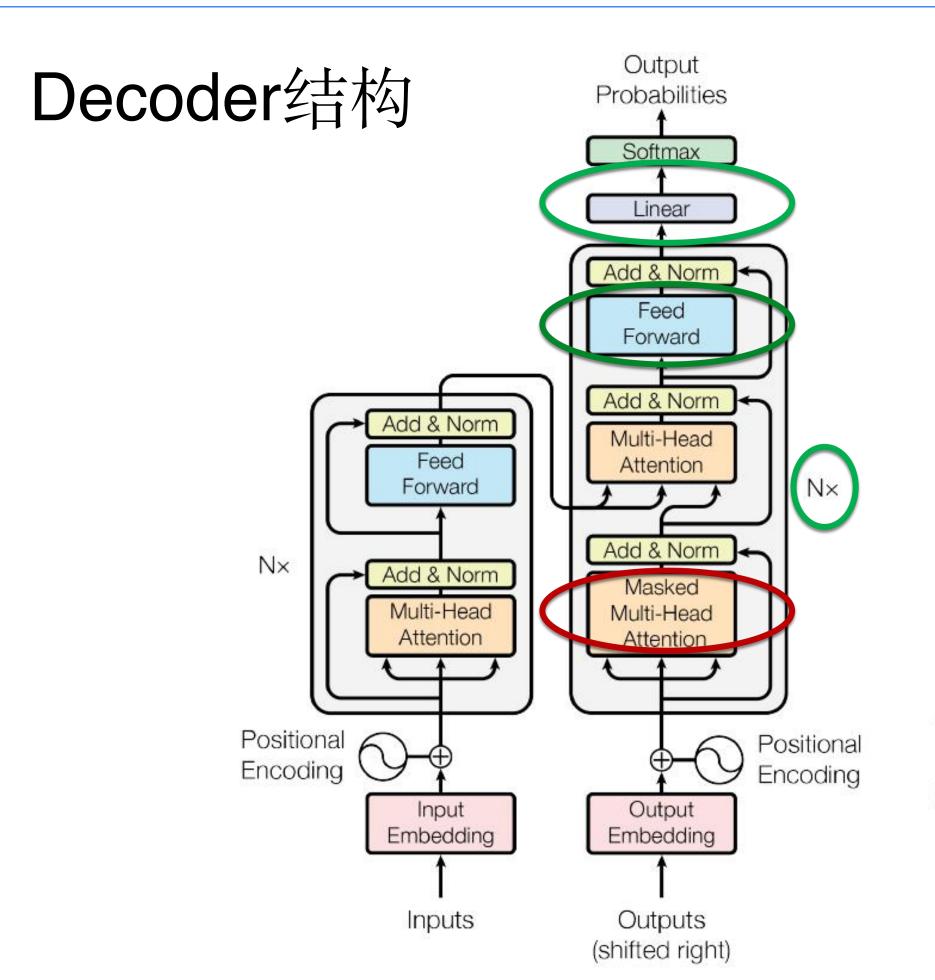
Yann N. Dauphin, Angela Fan, Michael Auli and David Grangier. "Language modeling with gated convolutional networks" International Conference on Machine Learning (2017). Biao Zhang, Deyi Xiong and Jinsong Su. "Accelerating Neural Transformer via an Average Attention Network" Meeting of the Association for Computational Linguistics (2018).

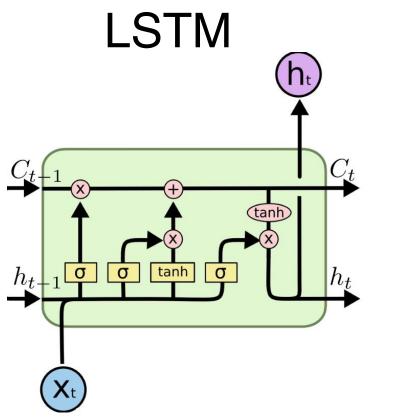


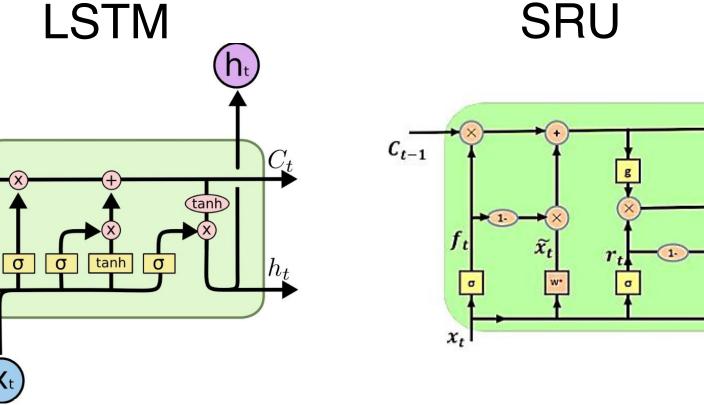


(shifted right)









$$f_{t} = \sigma_{g}(W_{f}x_{t} + U_{f}h_{t-1} + b_{f})$$

$$i_{t} = \sigma_{g}(W_{i}x_{t} + U_{i}h_{t-1} + b_{i})$$

$$o_{t} = \sigma_{g}(W_{o}x_{t} + U_{o}h_{t-1} + b_{o})$$

$$c_{t} = f_{t} \circ c_{t-1} + i_{t} \circ \sigma_{c}(W_{c}x_{t} + U_{c}h_{t-1} + b_{c})$$

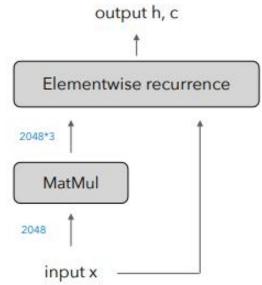
$$h_{t} = o_{t} \circ \sigma_{h}(c_{t})$$

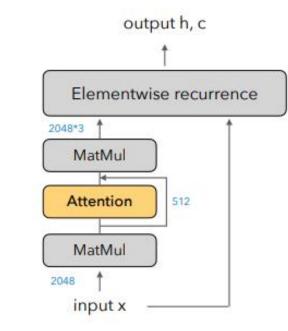
$$f_{t} = \sigma\left(\mathbf{W}_{f}x_{t} + \mathbf{v}_{f} \odot \mathbf{c}_{t-1} + \mathbf{b}_{f}\right)$$

$$c_{t} = f_{t} \odot \mathbf{c}_{t-1} + (1 - f_{t}) \odot (\mathbf{W}x_{t})$$

$$\mathbf{r}_{t} = \sigma\left(\mathbf{W}_{r}x_{t} + \mathbf{v}_{r} \odot \mathbf{c}_{t-1} + \mathbf{b}_{r}\right)$$

$$h_{t} = \mathbf{r}_{t} \odot \mathbf{c}_{t} + (1 - \mathbf{r}_{t}) \odot \mathbf{x}_{t}$$





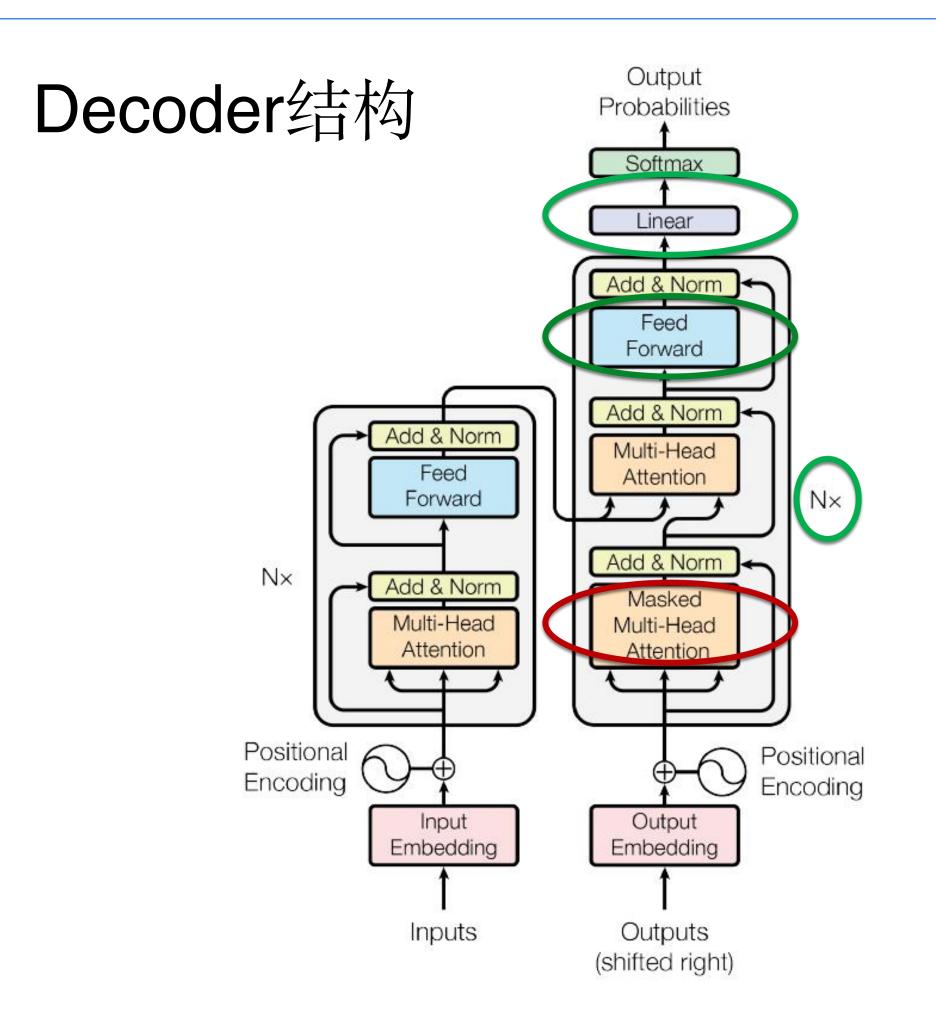
$$\mathbf{U}^ op = \left(egin{array}{c} \mathbf{W} \ \mathbf{W}_f \ \mathbf{W}_r \end{array}
ight) \left[\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_L
ight]$$

Tao Lei, Yu Zhang, Sida I. Wang, Hui Dai, Yoav Artzi, Simple Recurrent Units for Highly Parallelizable Recurrence, EMNLP 2017 Tao Lei. 2021. When Attention Meets Fast Recurrence: Training Language Models with Reduced Compute. Association for Computational Linguistics 2021.

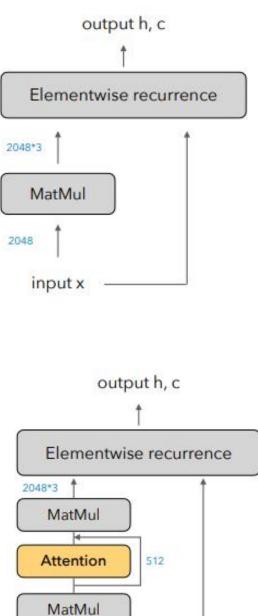








SRU++



AASRU

$$f[t] = \sigma(U[t, 0]) + V \odot C[t - 1] + b)$$

$$r[t] = \sigma(U[t, 1]) + V' \odot C[t - 1] + b')$$

$$c[t] = f[t] \odot C[t - 1] + (1 - f[t]) \odot C[t, 2]$$

$$h[t] = r[t] \odot C[t] + (1 - r[t]) \odot x[t]$$

$$Q = W^{q}X^{T}$$

$$V = W^{v}X^{T}$$

$$A^{T} = AVERAGE(V^{T})$$

$$U^{T} = W^{o}layernorm(Q + A)$$

Yann N. Dauphin, Angela Fan, Michael Auli and David Grangier. "Language modeling with gated convolutional networks" International Conference on Machine Learning (2017). Hengchao Shang, Ting Hu, Daimeng Wei, HW-TSC's Submission for the WMT22 Efficiency Task. In Proceedings of the Seventh Conference on Machine Translation (WMT), pages 677–681, ACL 2022

input x







策略	质量	速度	大小
知识蒸馏	$\sqrt{}$	$\sqrt{}$	_
量化推理	_	$\sqrt{}$	$\sqrt{}$
模型结构	$\sqrt{}$	$\sqrt{}$	_
参数共享	$\sqrt{}$	_	\checkmark
多语言	$\sqrt{}$	_	\checkmark
shortlist	_	\checkmark	_
Decoder结构	_	$\sqrt{}$	_



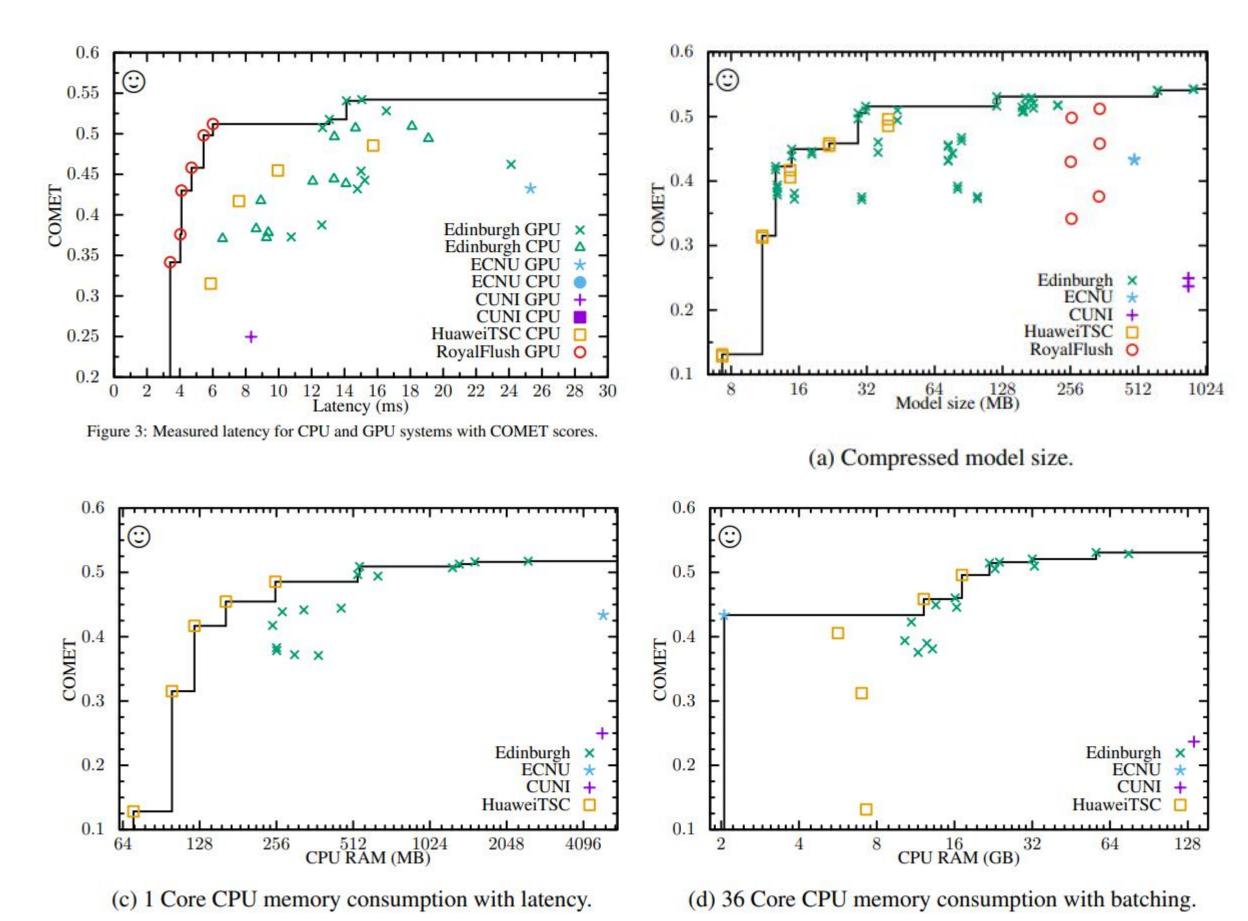




WMT22 Efficiency Task

Model	Precious	Size	WPS	BLEU
Teacher	FP32	2000	2	36.29
Base.12	FP32	212	237	35.30
	INT8	53	815	35.20
+retrain	INT8	53	815	35.19
Small.9	FP32	112	468	34.40
	INT8	28	1129	34.29
Tiny.12	FP32	68	759	33.62
	INT8	17	1693	33.39
Tiny.6	FP32	52	996	32.27
	INT8	13	2001	31.92
+int4	INT8	6.5	1989	26.10

Table 2: Optimization results. The test set is WMT 2020 News test. The unit of size is MB. WPS refers to the source side. The test environment is Intel(R) Xeon(R) Gold 6278C CPU @ 2.60GH. We submit four models: Base.12, Small.9, Tiny.12 and Tiny.6 and the final Tiny.6+int4



Hengchao Shang, Ting Hu, Daimeng Wei, Zongyao Li, Xianzhi Yu, Jianfei Feng, Ting Zhu, Lizhi Lei, Shimin Tao, Hao Yang, Ying Qin, Jinlong Yang, Zhiqiang Rao, and Zhengzhe Yu. 2022. HI Proceedings of the Seventh Conference on Machine Translation (WMT), pages 677–681, ACL 2022 Kenneth Heafield, Biao Zhang, Graeme Nail, Jelmer Van Der Linde, and Nikolay Bogoychev. 2022. Fficient Translation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 100–108, Abu Dhabi, United Arab Emirates (Hybrid). ACL 2022.







大纲

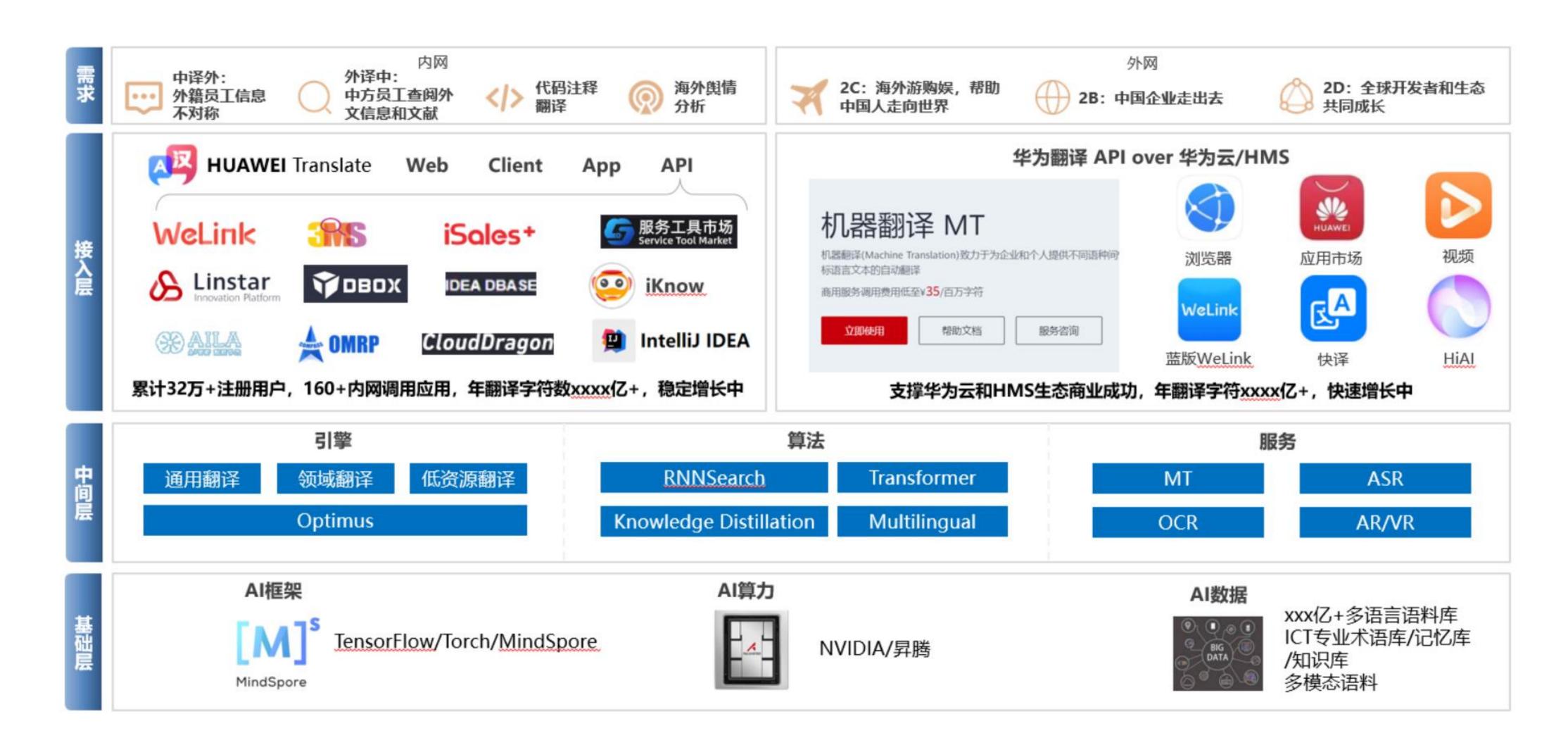
- 机器翻译简介
- 模型推理问题
- 端测推理加速
- 华为机器翻译
- 总结







华为机器翻译

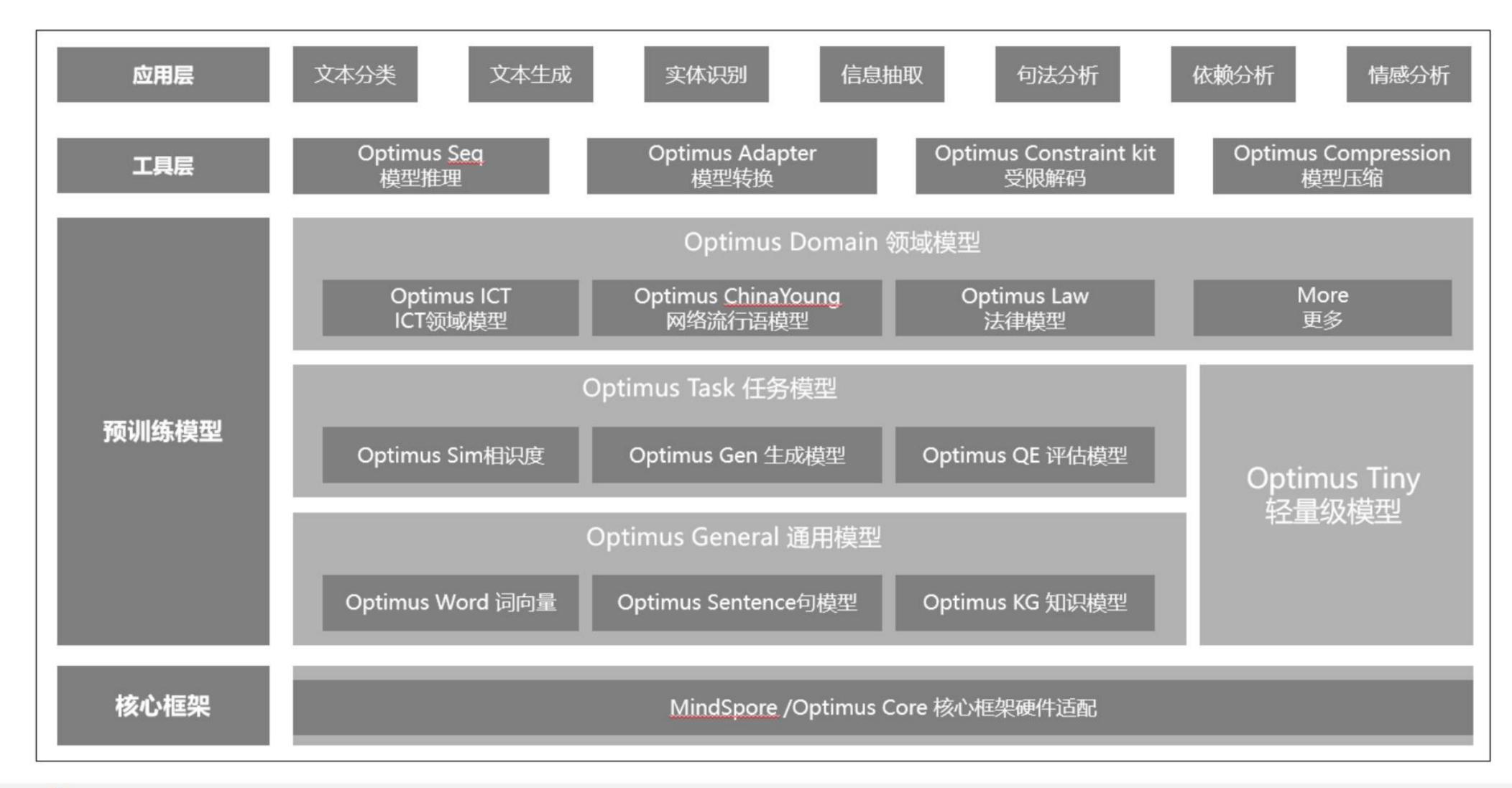








华为机器翻译









总结

- 1. 以业务为轴心, 按场景优化
- 2. 立足深度学习不断变化的大背景
- 3. 结合自身优势,不断迭代

策略	GPU	CPU	ARM
知识蒸馏	_	_	
量化推理	$\sqrt{}$	\checkmark	$\sqrt{}$
模型结构	$\sqrt{}$	$\sqrt{}$	$\sqrt{}$
参数共享	-	<u>-</u>	$\sqrt{}$
多语言	\checkmark	\checkmark	$\sqrt{}$
shortlist	_	\checkmark	\checkmark
Decoder结构	_	_	







想一想,我该如何把这些技术应用在工作实践中?

THANKS









【议题反馈】华为机器翻译模型训练推理加速实践

扫描二维码 提交议题反馈