

Trends of Machine Translation in China

—— NiuTrans for case study

Chunliang Zhang

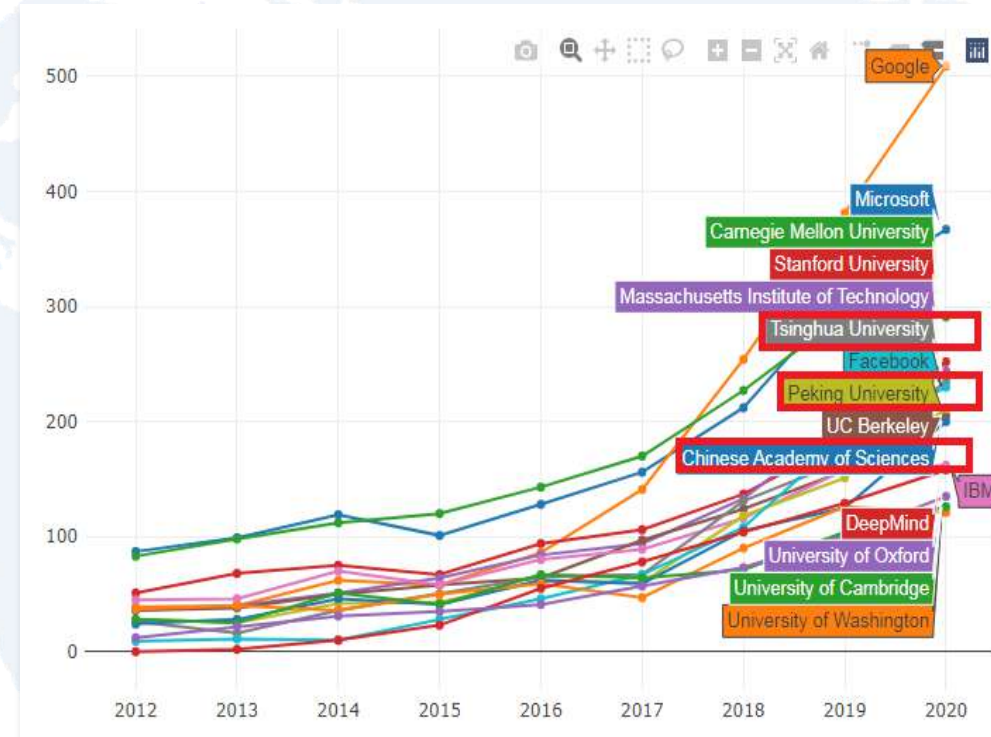
NiuTrans Co.
NLP Lab @ Northeastern University

1 NLP R&D power Comparison Worldwide

- By organization: 3 Chinese institutes among the top 15
- By country: The US taking the lead, followed by China and the UK

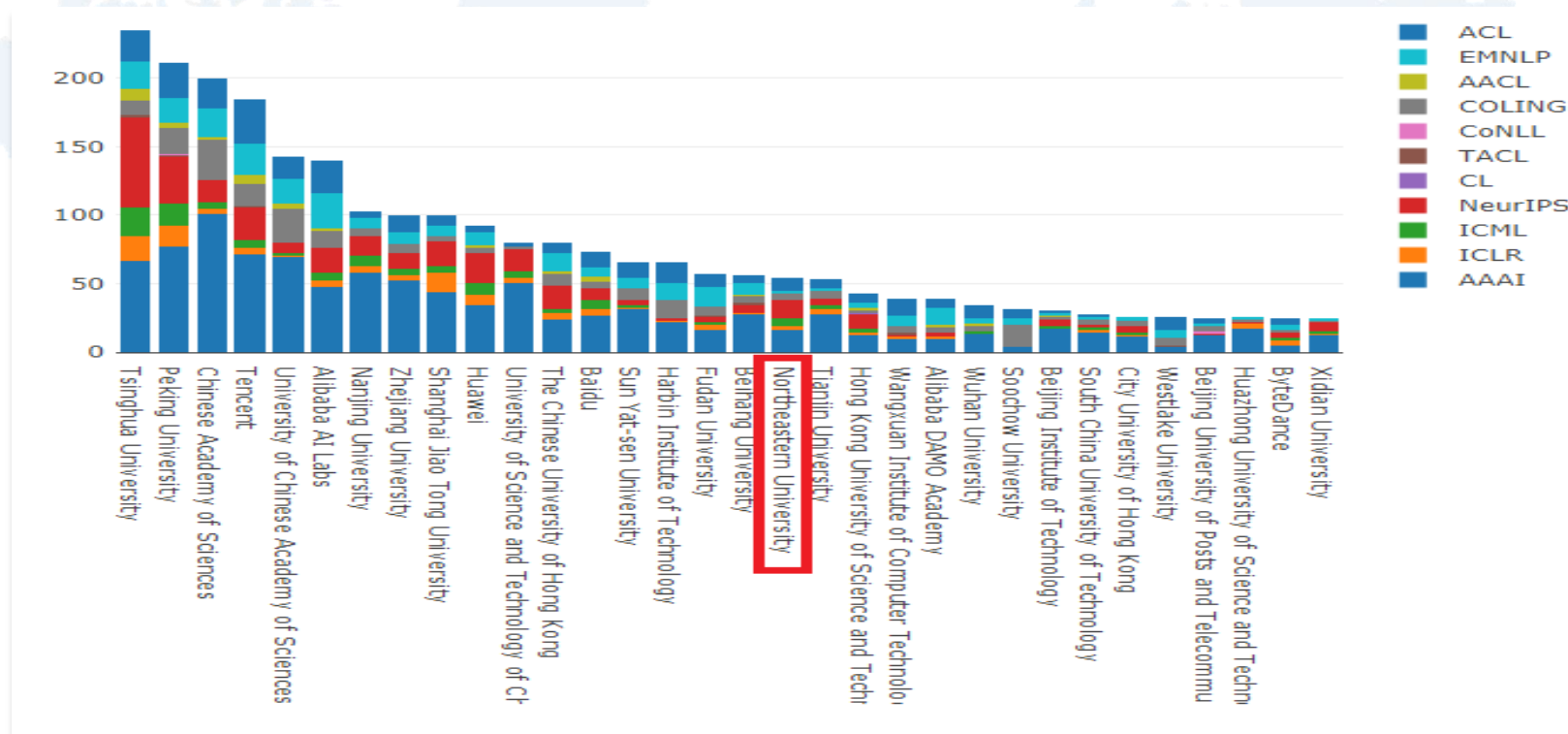


2012-2020 Paper publication by countries



2012-2020 Paper publication by organizations

- Most on the top 30 of the list are Universities and research institutes funded by government
- Only one university focus on Machine Translation: Northeastern University



The Top NLP Research Organizations in China (2012-2020)

<https://www.marekrei.com/blog/ml-and-nlp-publications-in-2020/>

3 Major MT Business Players in China

- Almost all the Internet-based service giants develop MT
- At least 3 teams are linked to universities/national institutes
- Some companies build hybrid business model: machine+human
- The focus is shifting from Tech R&D to MT engine, products & Marketing



Internet/AI Giants' Team



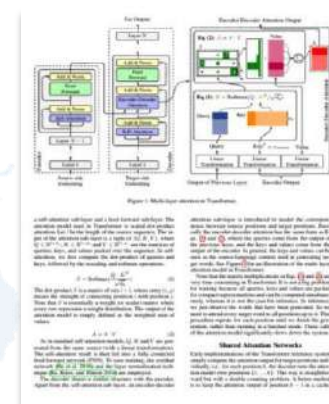
university-related Co.



the others (to name a few)

41 year for one thing: MT R&D

- NEUNLP lab was founded in 1980, a **NLP and MT research pioneer** in China
- A **leader of MT R&D into MT industry** in China
- 150+ members, 80% are PhDs and MSs
- 35 China central government research projects



304 languages

100+ major partners

- NENLP Lab founded

- built NiuTrans co.

- 22 languages, including 7 minority languages in China
- Write our own NiuTrans.NMT

- 118 languages
- NiuTrans Cloud Online
- Hangzhou, Shanghai, Shenzhen Offices

Dec. 1980

Oct. 2007

May 2012

2015

2016

2017

2018

2019-2020

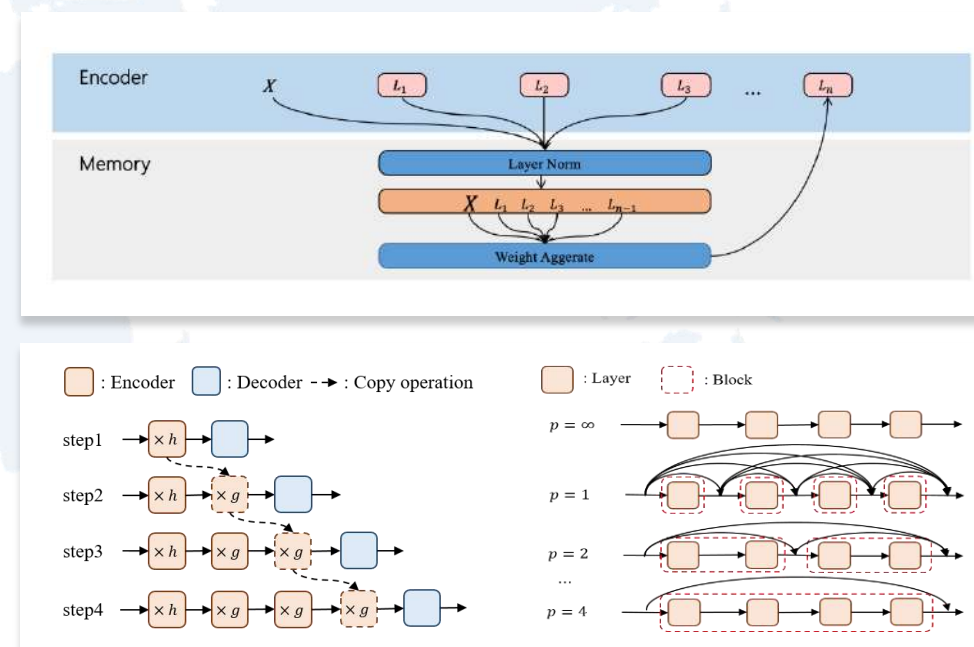
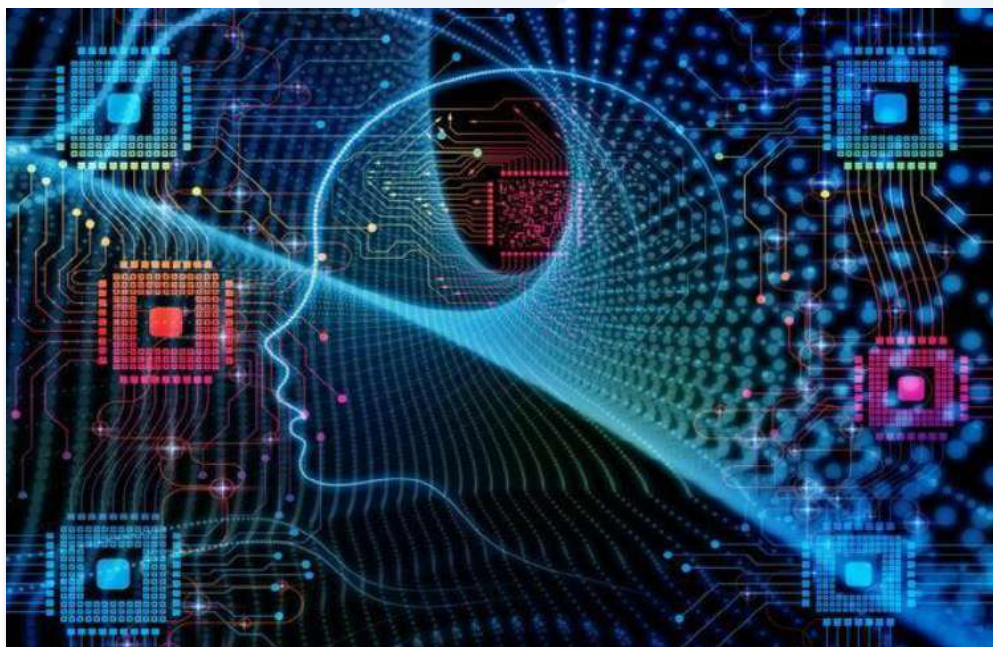
- NiuTrans.SMT @Lab

- 4 languages, CN, JN, KR, RU
- Angel investment by iFLYTEK

- 67 languages
- Beijing Office

- 304 languages
- Business partnership & go-international strategies

- **Deepening Transformer**
 - Achieve **better performance** by deepening the neural networks
 - Our model remains **stable** when it is especially deep (over 35 layers)
- **Efficient training methods for Deep Transformer**
 - Shallow-to-Deep Training: a simple yet **efficient** method



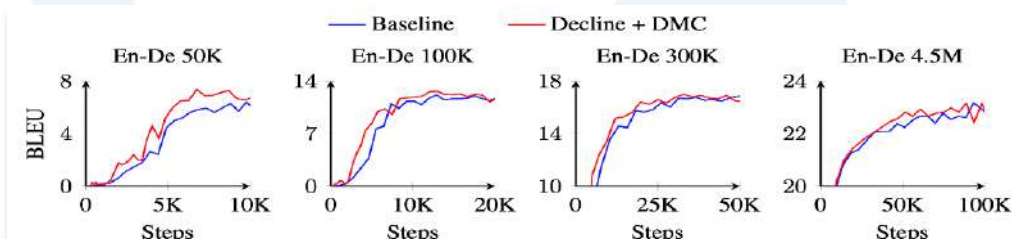
- **99% languages pairs do not have enough training data**
 - Corpus Building Technology: Bilingual Parallel corpus, Bilingual Comparable corpus, Monolingual corpus.
- **Bilingual Dictionary Induction and Dynamic Curriculum Learning**
 - Finding the corresponding words between two language by unsupervised method automatically.
 - Choosing the appropriate training samples for the current model.



Algorithm 1 Iterative Dimension Reduction

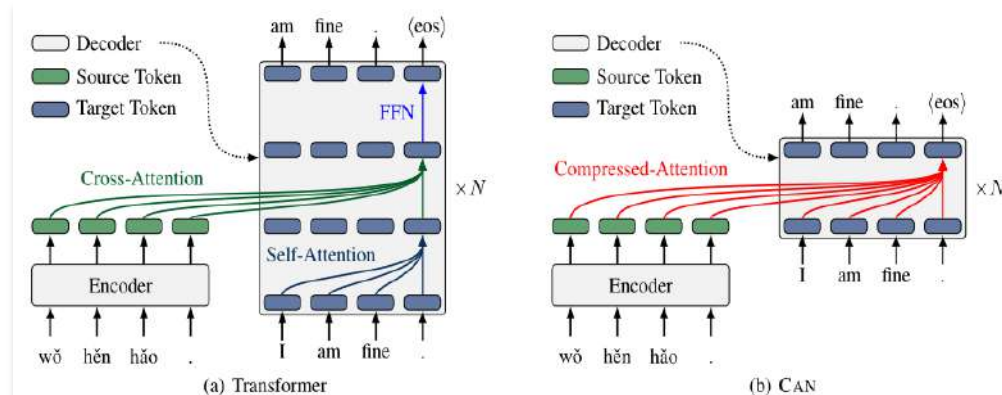
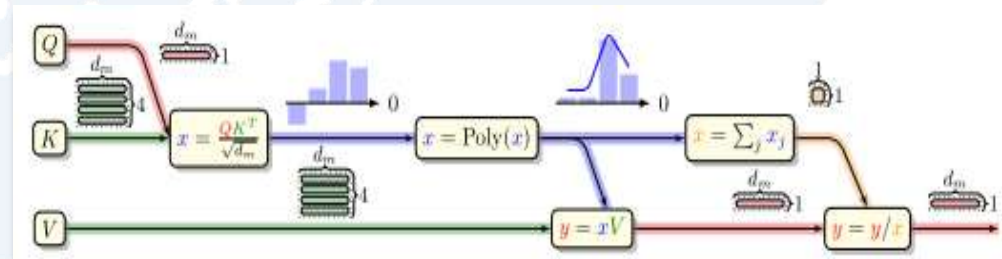
```

1: procedure IDR( $E, n$ )                                ▷  $E$  is the raw embeddings,  $n$  is the initial dimension
2:    $D \leftarrow \emptyset$                                 ▷ Set the dictionary to empty
3:   while  $n \leq 300$  do                                ▷ 300 is the dimension of the raw embeddings
4:     Reduce  $E$  to  $\bar{E}$  with dimension  $\min(n, 300)$  using PCA and dropout
5:     if  $D = \emptyset$  then
6:       Run the initialization and the self-learning on  $\bar{E}$ 
7:     else
8:       Run the self-learning on  $\bar{E}$  with  $D$  as the initial dictionary
9:     end if
10:    Translate 4K most frequent words and store the results in  $D$ 
11:     $n \leftarrow n \times 2$ 
12:  end while
13:  return  $W_X$  and  $W_Y$ 
14: end procedure
  
```



- Faster Transformer with less storage in decoding**

- Traditional Transformer relies on **massive computing power and storage (GPU)**.
- It's hard to apply the deep model in **mobile SoC (CPU)**.
- 8-bit Integer Inference for the Transformer Model + **lightweight model** architecture



MT Opensource, downloaded by 2000+ institutions in 70+ countries

- Started in 2007, NiuTrans.SMT v1.0 released in 2011
- The 1st successful Open source MT in China
- NiuTensor & NiuTrans.NMT** released in Dec. 2019

什么是张量

在计算机科学中，张量（Tensor）通常被定义为 n 维空间中的一种量，它具有 n 个分量，这种张量本质上是一个多维数组（array）。张量的阶或秩是这个多维数组的维度，或者简单理解为索引张量里的每个元素所需要的索引个数。通常来说，0阶张量（Scalar），1阶张量被定义为向量（vector），而2阶张量被定义为矩阵（matrix）。比如，在一个三维空间中，1阶张量的向量 (x, y, z) ，其中 x, y, z 分别表示这个点在三个轴上的坐标。

张量是一种高效的数学建模工具，它可以将复杂的问题通过统一、简洁的方式进行表达。比如，姜英俊同学做饭需要2斤的上牛肉每斤32元，土豆每斤2元，那么购买这些食物总共花费 $2 \times 32 + 5 \times 2 = 74$ 元。如果用张量来描述，我们可以用 $a = (2, 5)$ 表示所需不同食物的重量，然后用另一个1阶张量 $b = (32, 2)$ 表示不同食物的价格。最后，我们用一个0阶张量 c 表示总价，计算如下：

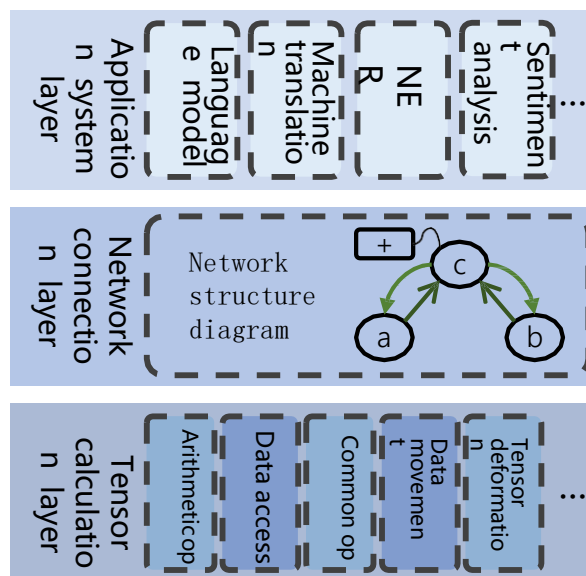
$$\begin{aligned} c &= a \times b^T \\ &= (2 \ 5) \times \begin{pmatrix} 32 \\ 2 \end{pmatrix} \\ &= 2 \times 32 + 5 \times 2 \\ &= 74 \end{aligned}$$

其中 b^T 表示行向量 b 的转置 - 列向量， \times 表示向量的乘法。第二天，姜英俊同学换了一个市场，这里牛肉每斤35元，土豆每斤1元。在两个市场分别购物的总价，可以把 a 重新定义为一个2阶张量 $\begin{pmatrix} 2 & 5 \\ 32 & 2 \end{pmatrix}$ ，总价 c 定义为一个2阶张量，同样有：

$$\begin{aligned} c &= a \times b^T \\ &= (2 \ 5) \times \begin{pmatrix} 32 & 35 \\ 2 & 1 \end{pmatrix} \\ &= (74 \ 75) \end{aligned}$$

即，在两个市场分别花费74元和75元。可以看出，利用张量可以对多样、复杂的问题进行建模，比如，可以进一步扩展上定义，把它们定义成更高阶的张量，处理不同时间、不同市场、不同菜品的情况，但是不论情况如何变化，都可以用同一套描述问题。

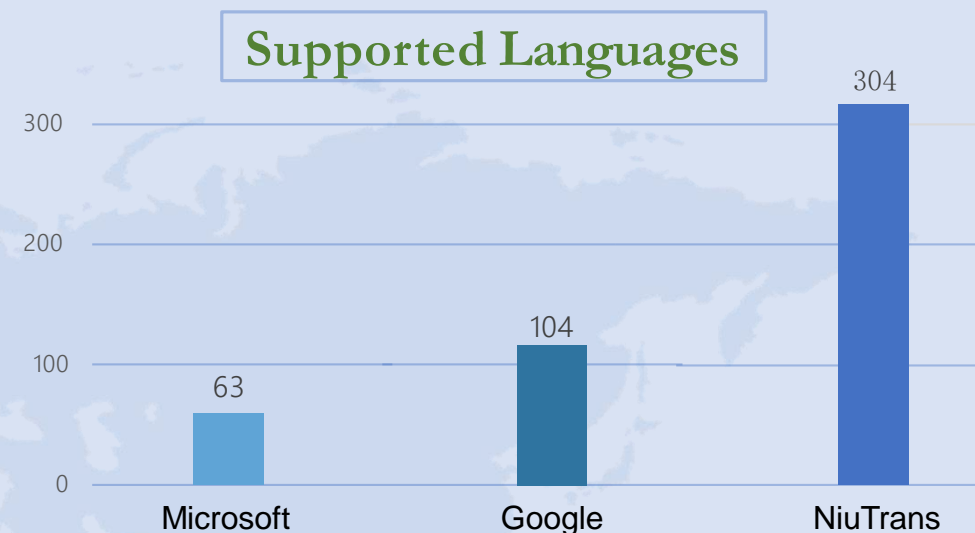
System architecture



Platform Comparison

	NiuTrans.NMT	OPEN-NMT (based on PyTorch)	Tensor2Tensor (based on TensorFlow)
Language	C++	Python	Python
Supported Device	CPU/GPU	CPU/GPU	CPU/GPU
Computation Graph	Dynamic	Dynamic	Static
Beam Search and Batch Decoding	✓	✓	✓
Extra ML Platform Dependence	X	✓	✓
GPU Memory for Decoding	1262MB	1300MB	1341MB

- Support **304** languages, translate between any two of them
- Most of them are **low-resource languages**, covering all China neighbor countries' official languages



China neighbors



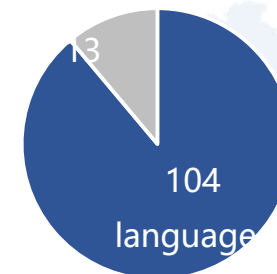
■ Covered

UN member countries



■ Covered

Language coverage in the official languages of the States members of the United Nations



Language covered

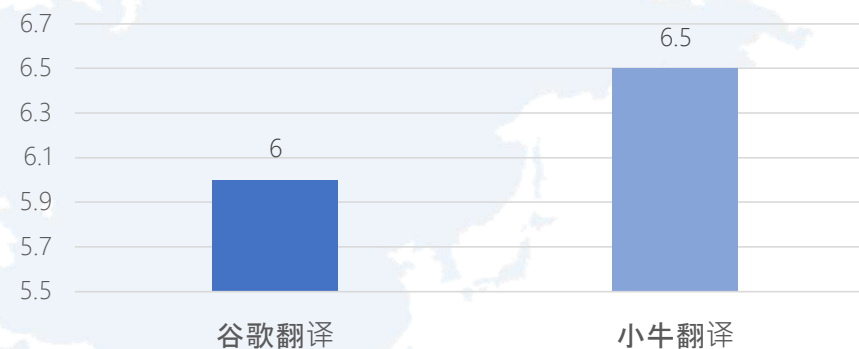
■ Covered ■ Uncovered

↖ Countries covered ↗

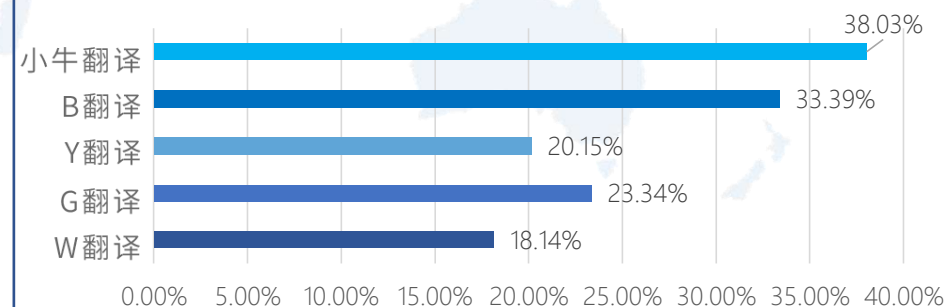
International
CompetitionWMT2018/2019/2020
NTCIR2011-JapanChina MT
Competition

CCMT2011/2018/2019

Human Evaluation 1



Human Evaluation 2





News Reports



Sports Commentary



A bite of local food



Hotel Service



Travel assistant

Medical service



Robot concierge



Easy transport



羽生結弦

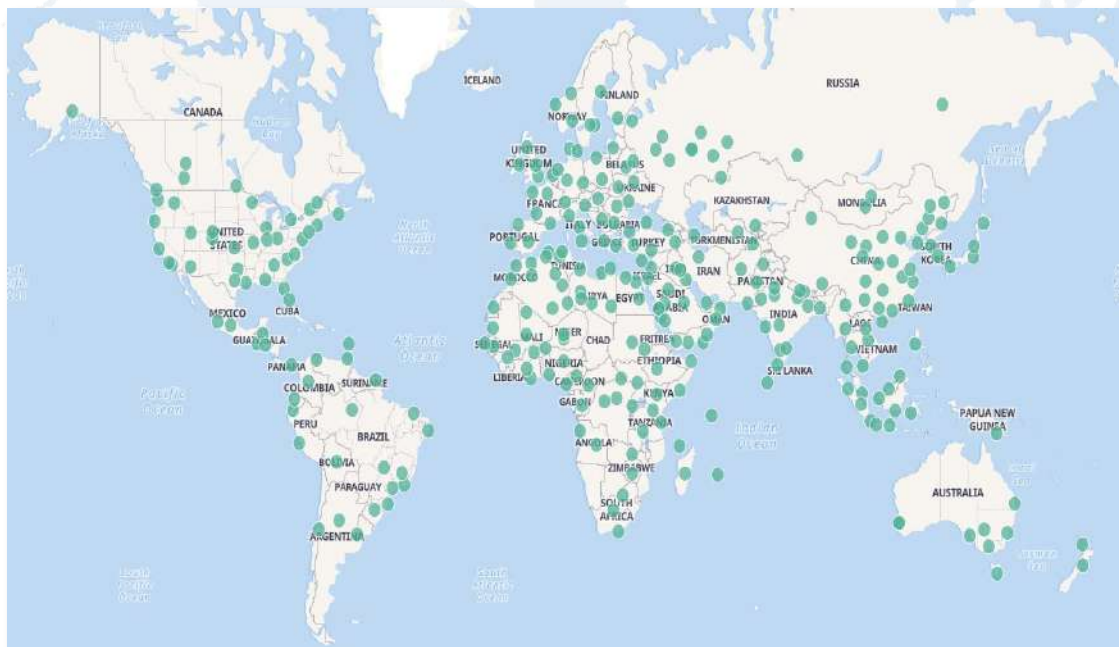
フィギュアスケート男子シングルの日本選手。2018年2月、平昌五輪フィギュアスケート男子シングルで優勝



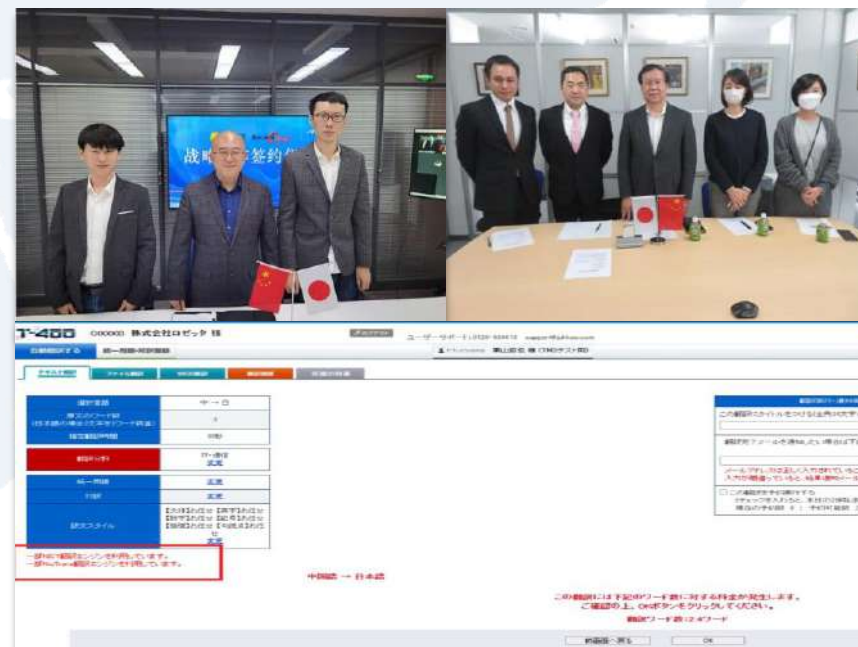
短道速滑

短道速滑は冬季奥运会项目，比赛采用淘汰制，以预、次、半决、决赛的比赛方式进行。

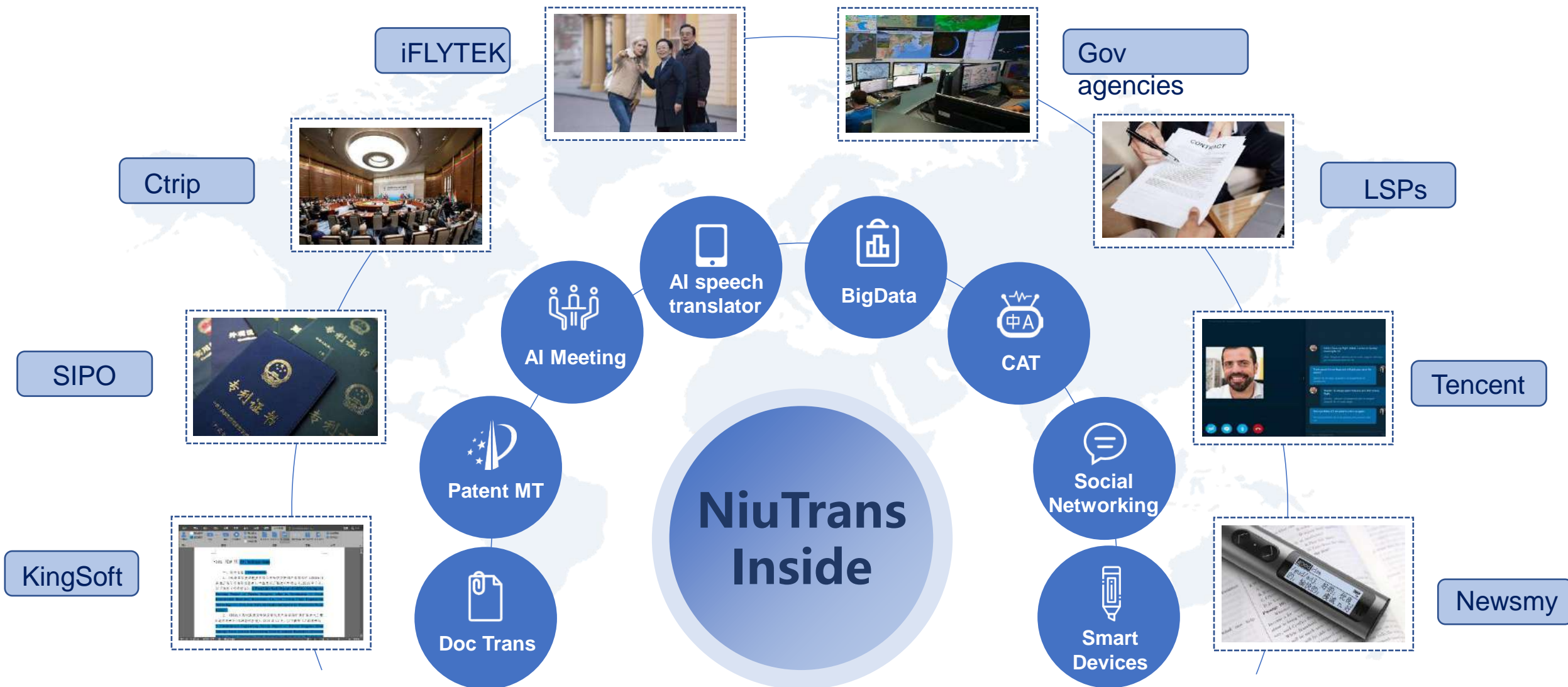
- 2016: Strategic Partnership with SYSTRAN Korea
- 2017: Release NiuTrans Plugin for Trados & memoQ
- 2020: Business Cooperation with 5 companies in Japan, including Rozetta



NiuTrans Cloud User Map



Partnership with Rozetta & RC Japan

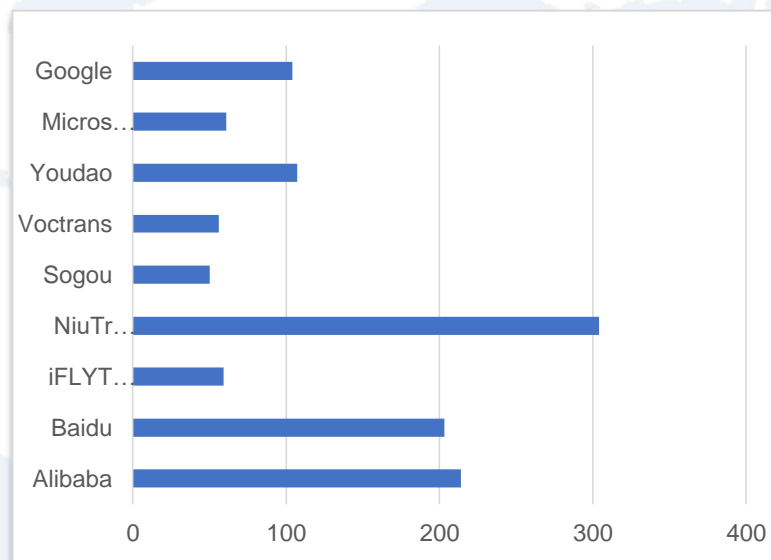


15 Post-pandemic MT Trends in China

- Much better MT quality is demanded for the low-resource languages
 - Hi-quality Corpus matter much more!



China's Belt&Road Map



MTs' Supported Language Statistics in 2020



High demand of language data

16 Post-pandemic MT Trends in China

- **More custom MT for certain businesses or industries**
 - Generic MT can't work the same good in all the fields



Medicine



Finance



Law

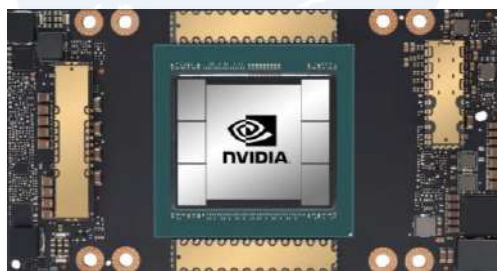


Engineering

17 Post-pandemic MT Trends in China

➤ AI Translation engine embedded into more devices

- integrate with other language tech: text + speech + image
- run on both GPU and CPU chips, especially CPU-powered smart mobile devices





小牛翻译
NiuTrans.com

Make the best MT engine

Thanks