



# MATHEMATICAL FRONTIERS

*The National  
Academies of* | SCIENCES  
ENGINEERING  
MEDICINE

[nas.edu/MathFrontiers](https://nas.edu/MathFrontiers)

Board on  
Mathematical Sciences & Analytics

# MATHEMATICAL FRONTIERS

## 2019 Monthly Webinar Series, 2-3pm ET

**February 12:** *Machine Learning for Materials Science\**

**March 12:** *Mathematics of Privacy\**

**April 9:** *Mathematics of Gravitational Waves\**

**May 14:** *Algebraic Geometry\**

**June 11:** *Mathematics of Transportation\**

**July 9:** *Cryptography & Cybersecurity\**

**August 13:** *Machine Learning in Medicine\**

**September 10:** *Logic and Foundations\**

**October 8:** *Mathematics of Quantum Physics\**

**November 12:** *Quantum Encryption\**

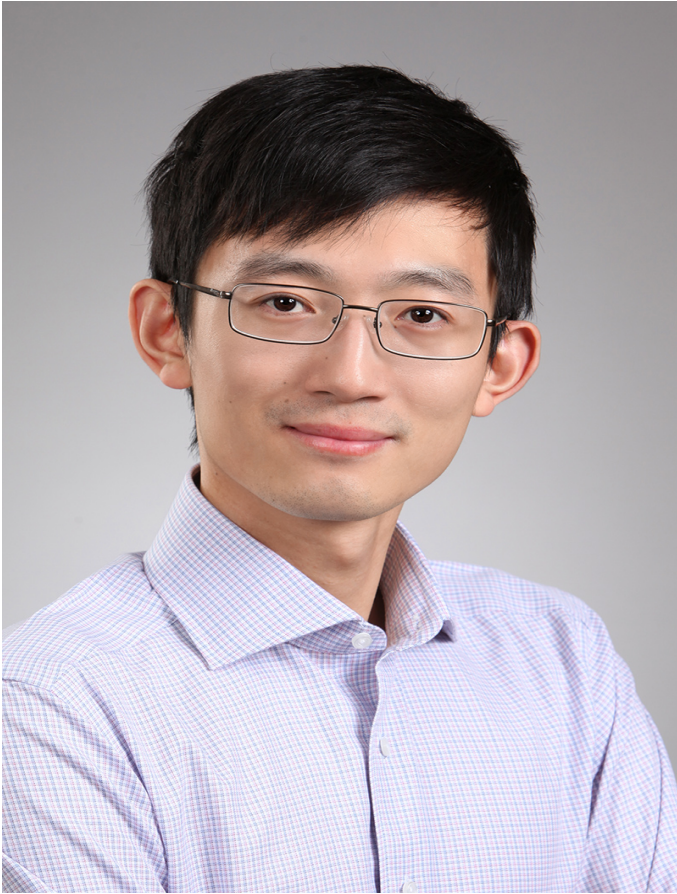
**December 10:** *Machine Learning for Text*

*Made possible by support for BMSA from the  
**National Science Foundation**  
**Division of Mathematical Sciences**  
and the  
**Department of Energy**  
**Advanced Scientific Computing Research***

*\* Webinar posted*

# MATHEMATICAL FRONTIERS

## Machine Learning for Text



**Tengyu Ma,  
Stanford University**



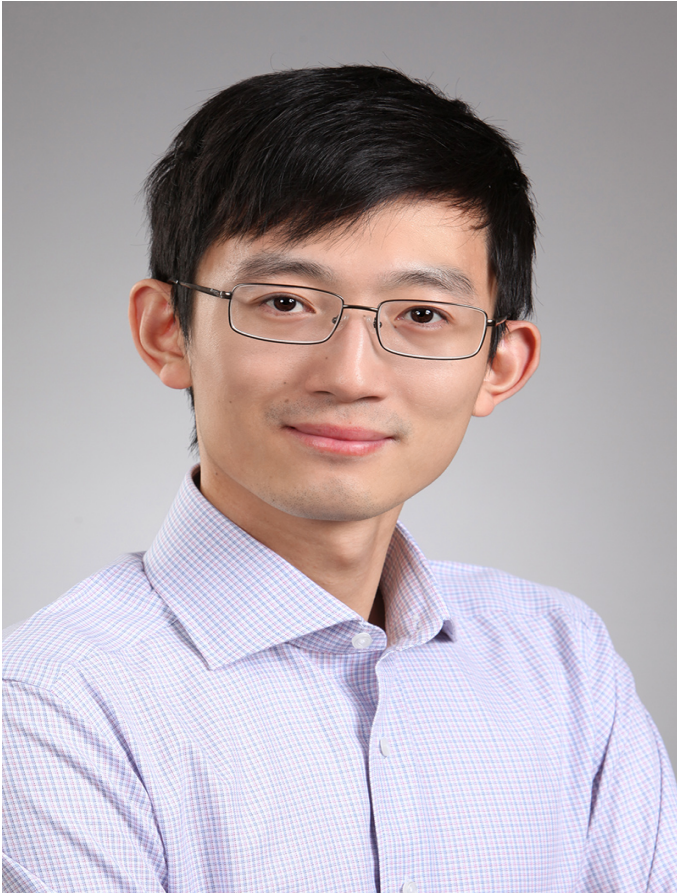
**Marine Carpuat,  
University of Maryland**



**Mark Green,  
UCLA (moderator)**

# MATHEMATICAL FRONTIERS

## Machine Learning for Text



**Tengyu Ma,**  
**Stanford University**

*Assistant Professor of Computer Science and Statistics*

## **Machine Learning for Texts: Understanding Embeddings**

# Breakthroughs in Natural Language Processing

Google Translate

DETECT LANGUAGE

FRENCH

CHINESE

ENGLISH



CHINESE (SIMPLIFIED)

FRENCH

ENGLISH



This is a talk about machine learning for texts.



48/5000



Ceci est une discussion sur l'apprentissage automatique pour les textes.



ENGLISH - DETECTED

FRENCH

CHINESE

ENGLISH



CHINESE (SIMPLIFIED)

FRENCH

ENGLISH



This is a talk about machine learning for texts|



47/5000



这是关于文本机器学习的话题

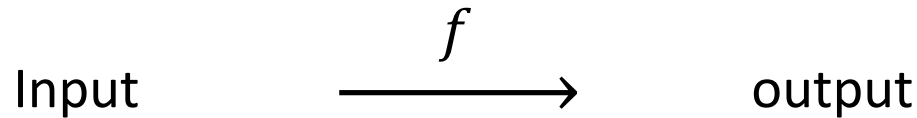


Zhè shì guānyú wénběn jīqì xuéxí de huàtí

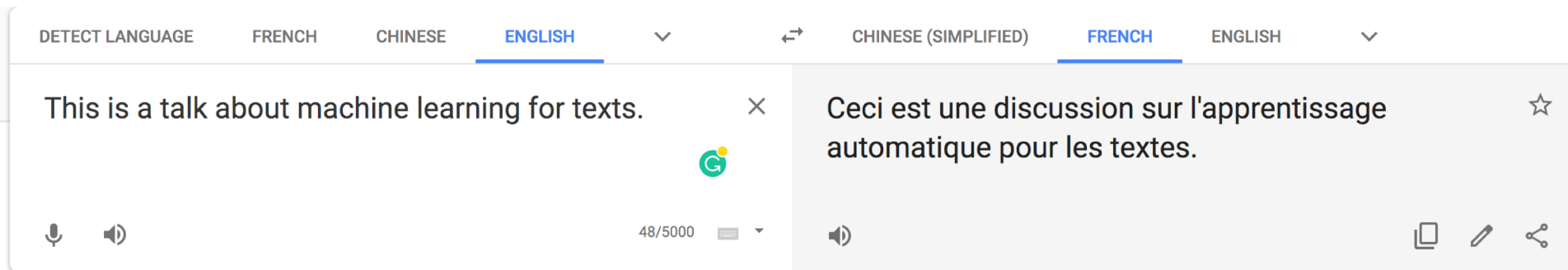


# Machine Learning (Supervised Learning)

- Find a function  $f$



## ➤ Translation



## ➤ Sentiment analysis

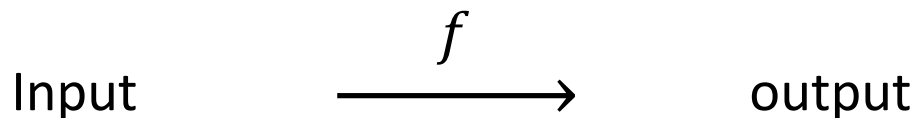
This seminar series is fantastic.



positive  
sentiment

# Machine Learning (Supervised Learning)

➤ Find a function  $f$



How do we represent **texts inputs** as numerical values?

# Classic “One-hot” and “Bag-of-words” Representation

- Vocabulary = {a, aardvark, aardwolf, ..., zymurgy} of size  $N$

$$\text{happy} \longrightarrow v_{\text{happy}} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ \vdots \\ 1 \\ \vdots \\ \vdots \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} \text{a} \\ \text{aardvark} \\ \text{aardwolf} \\ \\ \\ \text{happy} \\ \\ \\ \text{zymurgy} \end{matrix} \in \mathbb{R}^N$$

# Classic “One-hot” and “Bag-of-words” Representation

- Vocabulary = {a, aardvark, aardwolf, ..., zymurgy} of size  $N$

A happy  
aardwolf

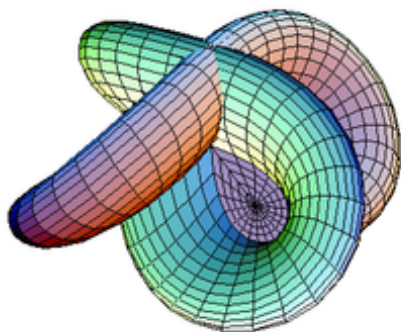
→  $v =$

$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \\ \vdots \\ 1 \\ \vdots \\ \vdots \\ \vdots \\ 0 \end{bmatrix}$$

a  
aardvark  
aardwolf  
  
happy  
  
zymurgy

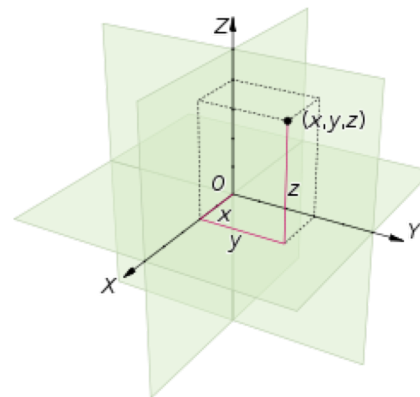
$\in \mathbb{R}^N$

# Embeddings (in Machine Learning)



$$x \in \mathcal{X}$$

complicated space



$$v_x \in \mathbb{R}^d$$

Euclidean space with  
**meaningful** inner products

# Word Embeddings

Vocabulary



$\mathbb{R}^{300}$

Goal:

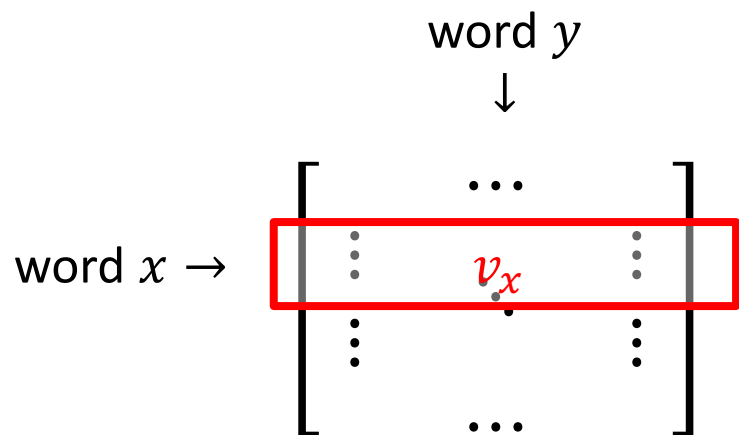
embedding captures semantics information  
(ideally via linear algebraic operations)

- inner products characterize similarity
  - similar words have large inner products
- differences characterize relationship
  - analogous pairs have similar differences



# Distributional Hypothesis of Meaning ([Harris'54], [Firth'57])

Meaning of a word is determined by words it **co-occurs** with.



Co-occurrence matrix  
 $\text{Pr}(\cdot, \cdot)$

**Def:**  $\text{Pr}(x, y) \triangleq$  prob. of co-occurrences  
of  $x, y$  in a window of size 5

“a window of size 5”

➤ Rows of co-occurrence matrix are  
reasonable embeddings [Lund-Burgess'96]

➤ [Church-Hanks'90]

$$v_x = \text{row of PMI}(x, y) \triangleq \log \frac{\text{Pr}[x, y]}{\text{Pr}[x] \text{Pr}[y]}$$

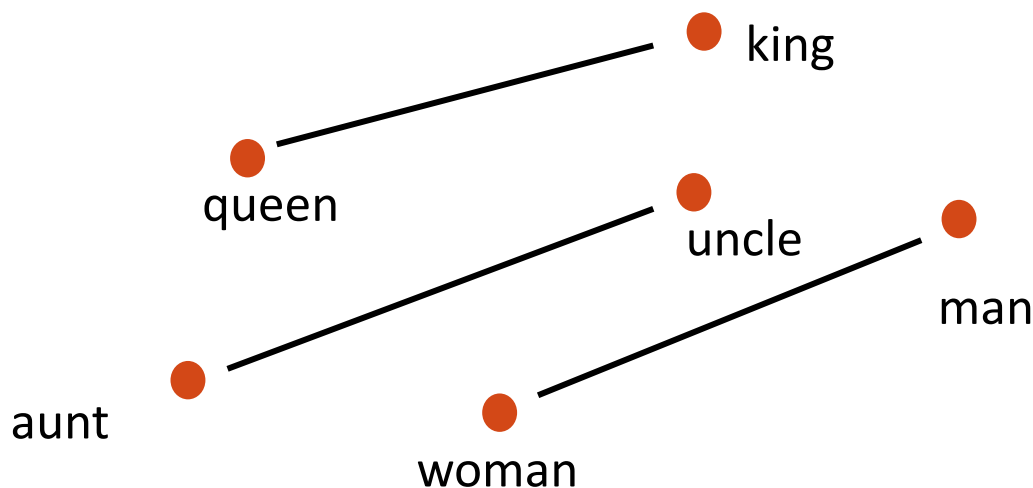
(PMI = point-wise mutual information)

# Dimension-Reduced PMI Embeddings [Levy-Goldberg'14]

1. Compute  $\text{PMI}(x, y) = \log \frac{\text{Pr}[x, y]}{\text{Pr}[x] \text{Pr}[y]}$
2. Take rank-300 SVD (best rank-300 approximation) of PMI
  - $\Leftrightarrow$  Fit  $\text{PMI}(x, y) \approx \langle v_x, v_y \rangle$  where  $v_x \in \mathbb{R}^{300}$

➤ “Linear structure” in the found  $v_x$ 's :

$$v_{\text{woman}} - v_{\text{man}} \approx v_{\text{queen}} - v_{\text{king}} \approx v_{\text{uncle}} - v_{\text{aunt}} \approx \dots$$



# Non-linear Embedding methods

- word2vec [Mikolov et al'13] :

$$\Pr[ x_{i+6} \mid x_{i+1}, \dots, x_{i+5} ] \propto \exp \langle v_{x_{i+6}}, \frac{1}{5} (v_{x_{i+1}} + \dots + v_{x_{i+5}}) \rangle$$

- GloVe [Pennington et al'14] :

$$\log \Pr[x, y] \approx \langle v_x, v_y \rangle + s_x + s_y + C$$

- [Levy-Goldberg'14] (Previous slide)

$$\text{PMI}(x, y) = \log \frac{\Pr[x, y]}{\Pr[x] \Pr[y]} \approx \langle v_x, v_y \rangle + C$$

Logarithm (or exponential) seems to exclude linear algebra!

# Where does the log come from?

[Arora et al.'16, c.f. Levy-Goldberg'14, Pennington et al'14]

➤ For most of the words  $\chi$ :

$$\frac{\Pr[\chi \mid \text{king}]}{\Pr[\chi \mid \text{queen}]} \approx \frac{\Pr[\chi \mid \text{man}]}{\Pr[\chi \mid \text{woman}]}$$

➤ For  $\chi$  unrelated to gender: LHS, RHS  $\approx 1$

➤ for  $\chi = \text{dress}$ , LHS, RHS  $\ll 1$  ; for  $\chi = \text{John}$ , LHS, RHS  $\gg 1$

$$\Rightarrow \sum_{\chi} \left( \log \frac{\Pr[\chi \mid \text{king}]}{\Pr[\chi \mid \text{queen}]} - \log \frac{\Pr[\chi \mid \text{man}]}{\Pr[\chi \mid \text{woman}]} \right)^2$$

$$\| \text{PMI}(\text{king}, \cdot) - \text{PMI}(\text{queen}, \cdot) - \text{PMI}(\text{man}, \cdot) + \text{PMI}(\text{woman}, \cdot) \|_2^2 \approx 0$$

➤ Rows of PMI matrix has “linear structure”

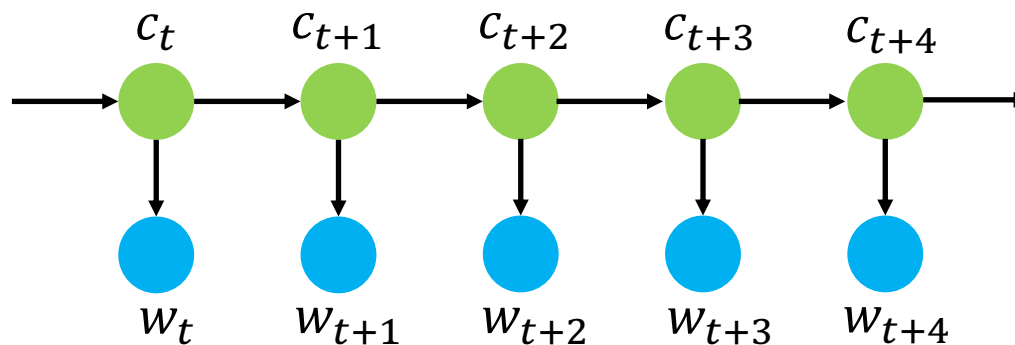
# Why Does Dimension Reduction Help?

Empirically can find vectors  $v_x$ 's such that

$$\text{PMI}(x, y) \approx \langle v_x, v_y \rangle$$

1. PMI is not necessarily PSD
2. Relative approximation error is high (17%); Low-dimensional  $v_x$ 's have **better** linear structure than rows of PMI

# RAND-WALK: A Generative Model for Language [Arora et al'16]



## ➤ Hidden Markov Model:

- discourse vector  $c_t \in \mathbb{R}^d$  governs the discourse/theme/context of time  $t$
- words  $w_t$  (observable); embedding  $v_{w_t} \in \mathbb{R}^d$  (parameters to learn)
- log-linear observation model

$$\Pr[w_t \mid c_t] \propto \exp\langle v_{w_t}, c_t \rangle$$

Closely related to [Mnih-Hinton'07]

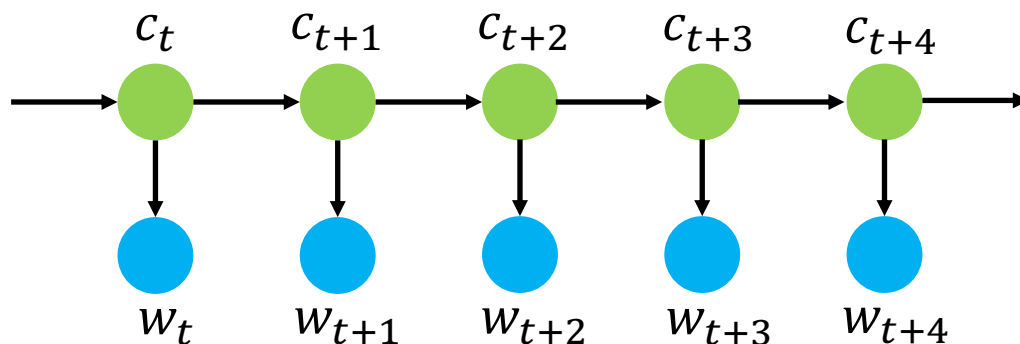
# Why Does Dimension Reduction Help?

Empirically can find vectors  $v_x$ 's such that

$$\text{PMI}(x, y) \approx \langle v_x, v_y \rangle$$

1. PMI is not necessarily PSD and low-rank
  - Under rand-walk model, PMI is approximately PSD and low-rank
2. Relative approximation error is high (17%); Low-dimensional  $v_x$ 's have better linear structure than rows of PMI
  - Dimension-reduction reduces the noises

**RAND-WALK  
Model**



# Summary and Looking Ahead

- Theoretical explanations of embeddings methods
  - Popular embeddings methods, such as PMI+SVD, word2vec, Glove can be viewed as algorithms for learning a generative model of language
- Follow-up works: embeddings for sentences, polysemous words, rare words [Arora et al.'17,18a&b ...]
- Open directions:
  - Understanding the state-of-the-art contextualized embeddings (Elmo, Bert, etc..)
  - Optimizations of the embeddings
  - Understanding other algorithms for other tasks in NLP (machine translation, etc.)
  - A theory of representation learning

# Main References

- RAND-WALK: A Latent Variable Model Approach to Word Embeddings. Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. Transactions of the Association for Computational Linguistics (TACL), 2016
- A Simple but Tough-to-Beat Baseline for Sentence Embeddings. Sanjeev Arora, Yingyu Liang, Tengyu Ma. International Conference on Learning Representations (ICLR) 2017
- Linear Algebraic Structure of Word Senses, with Applications to Polysemy. Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. TACL, 2018
- A La Carte Embedding: Cheap but Effective Induction of Semantic Feature Vectors. Mikhail Khodak, Nikunj Saunshi, Yingyu Liang, Tengyu Ma, Brandon Stewart, Sanjeev Arora. ACL, 2018
- Neural Word Embedding as Implicit Matrix Factorization. Omer Levy and Yoav Goldberg. Neurips 2014.
- GloVe: Global Vectors for Word Representation. Jeffrey Pennington, Richard Socher, and Christopher D. Manning. EMNLP, 2014.
- Distributed Representations of Words and Phrases and their Compositionality. Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, Jeffrey Dean. NIPS 2013.

# MATHEMATICAL FRONTIERS

## Machine Learning for Text



Marine Carpuat,  
University of Maryland

*Assistant Professor in Computer Science*  
*marine@cs.umd.edu*

## **Toward Breaking Language Barriers with Neural Machine Translation**



Search for a language, dialect name or major city...



这座中国首都拥有速度高得惊人的互联网，有人脸识别软件等尖端技术，在人工智能方面投入了巨资并且拥有无可匹敌的国际化能量，它对富于探索精神的外国人而言是最激动人心的城市之一。

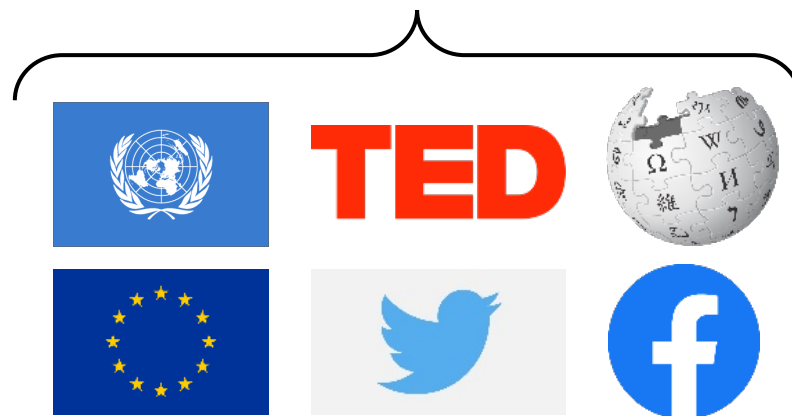
The Chinese capital, with its surprisingly high-speed Internet, sophisticated technology such as face-recognition software, has invested heavily in artificial intelligence and has unrivaled international energy, and is one of the most exciting cities for exploration-minded foreigners.

# Translation as Machine Learning

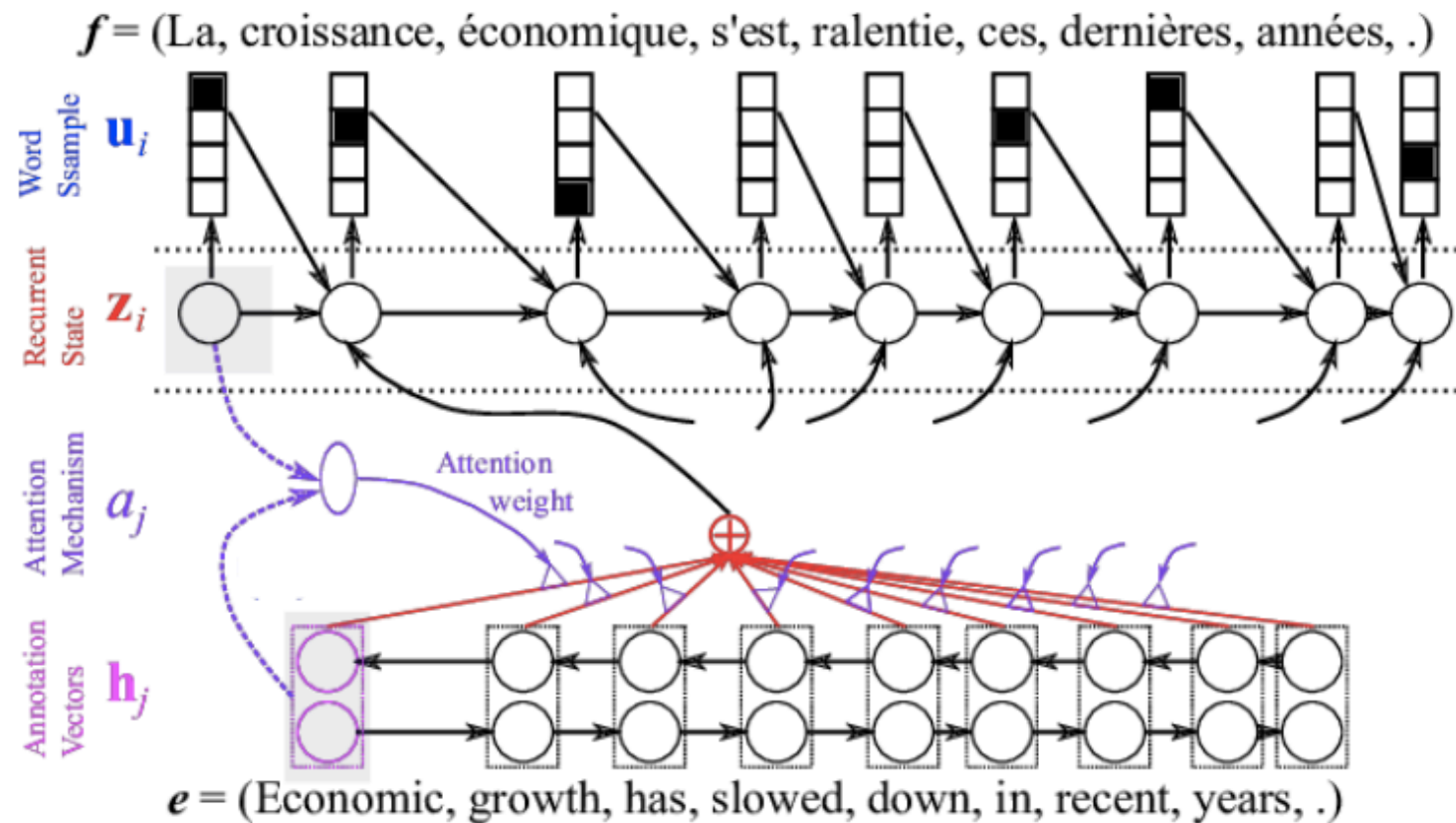
An English sentence  $e$  is translated  
into the French sentence

$$f^* = \operatorname{argmax}_f p(f|e; \theta)$$

$$\theta^* = \operatorname{argmax}_\theta \sum_i \log p(f_i | e_i; \theta)$$



# Translation as Deep Learning



$$p(f | e; \theta) = \prod_{t=1}^{|f|} p(f_t | f_{<t}, e; \theta)$$

# Translation as Deep Learning: Challenges

requires millions of translation  
examples

not available for many languages!

raises fundamental machine learning  
challenges

intractably large output space, infinitely  
many correct outputs...

makes errors that have real world  
impact

yet models are opaque, and developed  
independency from use cases

Toward better  
translation with  
limited training  
data

Some approaches:

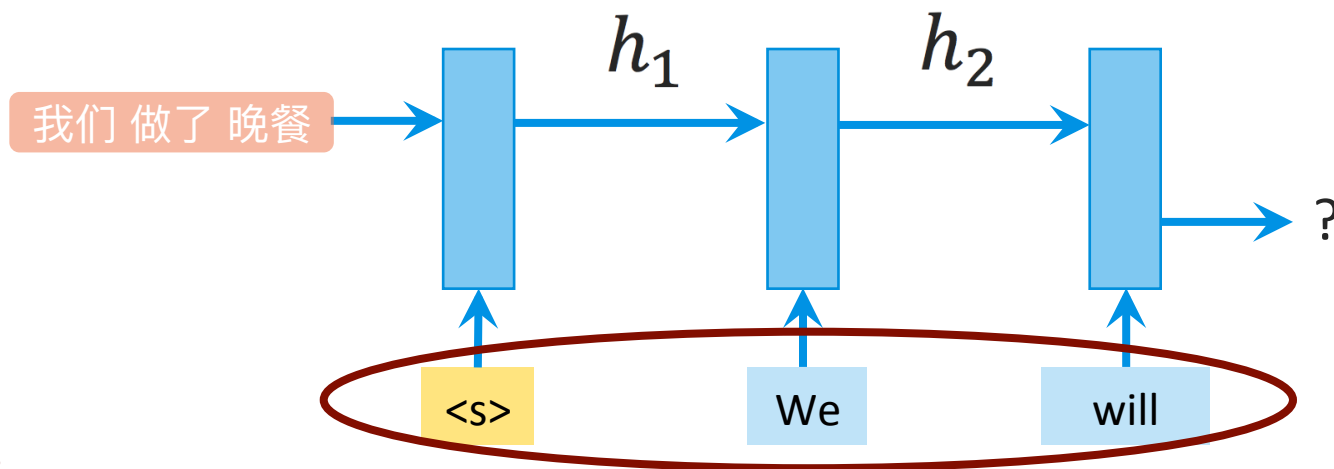
Learn from related languages

Learn from monolingual text

**Improve the training objective**

# Training Problem: Exposure Bias, a Gap Between Training and Inference

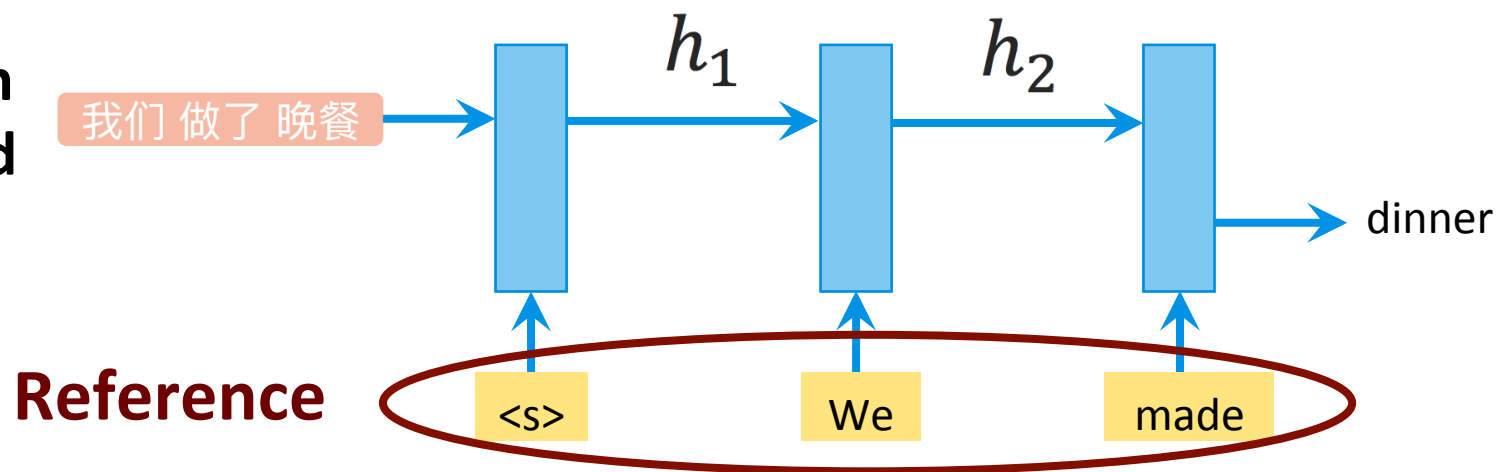
Inference



**Model  
Translation**

$$P(f|e) = \prod_{t=1}^T p(f_t | f_{<t}, e)$$

Maximum  
Likelihood  
Training



**Reference**

**Loss =**

$$\sum_{t=1}^T \log p(f_t | f_{<t}, e)$$

# How to Address Exposure Bias?

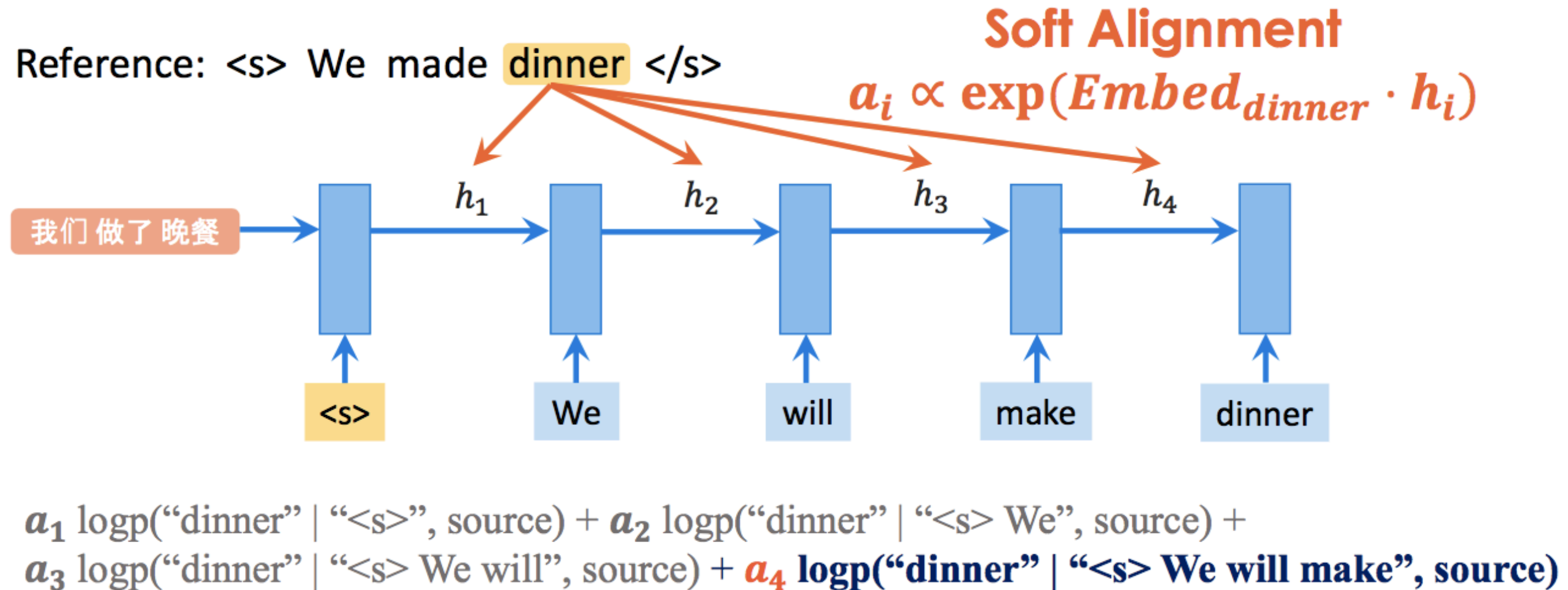
Expose models to their own predictions during training

But how to compute the loss when the partial translation diverges from the reference?

Our method:

1. **Generate translation prefixes** via differentiable sampling
2. Learn to **align** the reference words with sampled prefixes

# Our Solution: Align Reference with Partial Translations



Toward better  
translation with  
limited training  
data

Some approaches:

Learn from related languages

Learn from monolingual text

**Improve the training objective**

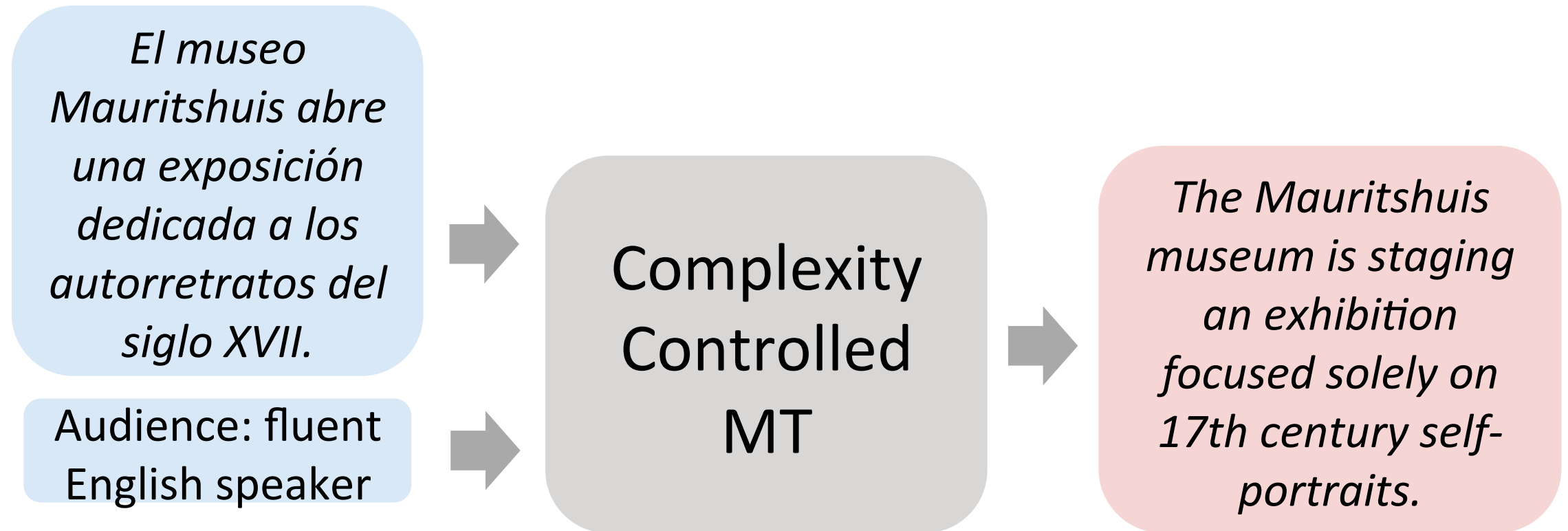
# Toward more user-centered machine translation

Can machine translation help  
human translators and interpreters  
be more productive?

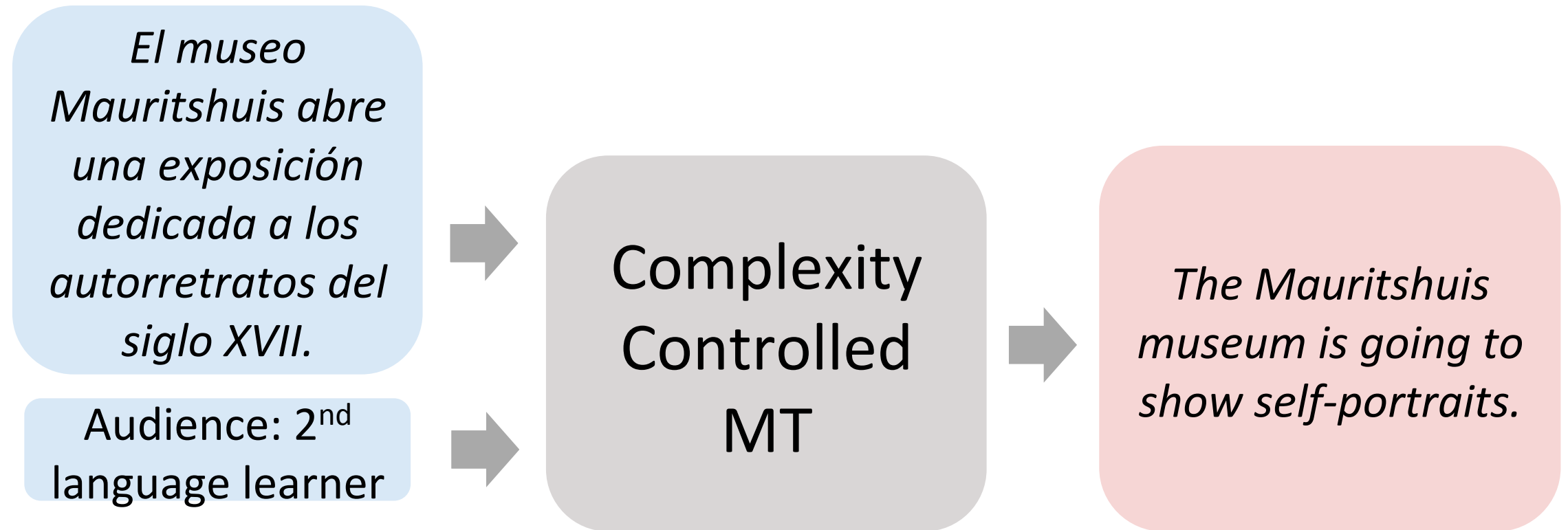
What errors matter most for  
different use cases?

Can we tailor machine translation  
output to different audiences?

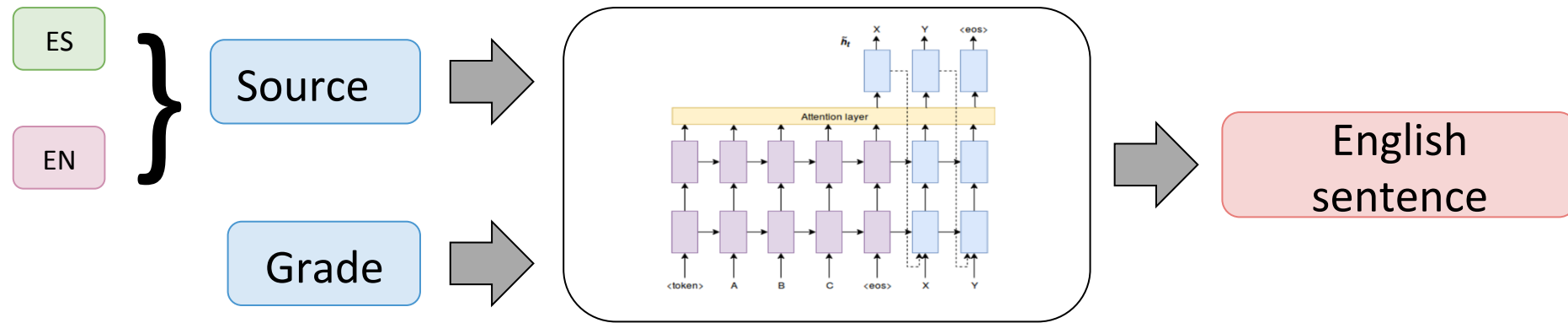
# Controlling MT Complexity for Different Audiences



# Controlling MT Complexity for Different Audiences



# Adapting translation output to different audiences via multi-task learning



$$\text{Multi-task loss} = \underbrace{\sum_{(s_i, g_e, e_o)} \log P(e_o | s_i, g_e; \theta)}_{L_{CMT}} + \underbrace{\sum_{(e_i, g_e, e_o)} \log P(e_o | e_i, g_e; \theta)}_{L_{Simplify}} + \underbrace{\sum_{(s_i, e_o)} \log P(e_o | s_i; \theta)}_{L_{MT}}$$

Spanish sentences translated into simpler English

Complex English sentences paired with simpler English

Spanish-English translation examples



# References

Sweta Agrawal and Marine Carpuat.

[“Controlling Text Complexity in Neural Machine Translation”](#). EMNLP 2019.

Weijia Xu, Xing Niu and Marine Carpuat.

[“Differentiable Sampling with Flexible Reference Word Order for Neural Machine Translation”](#). NAACL 2019.

Weijia Xu and Marine Carpuat.

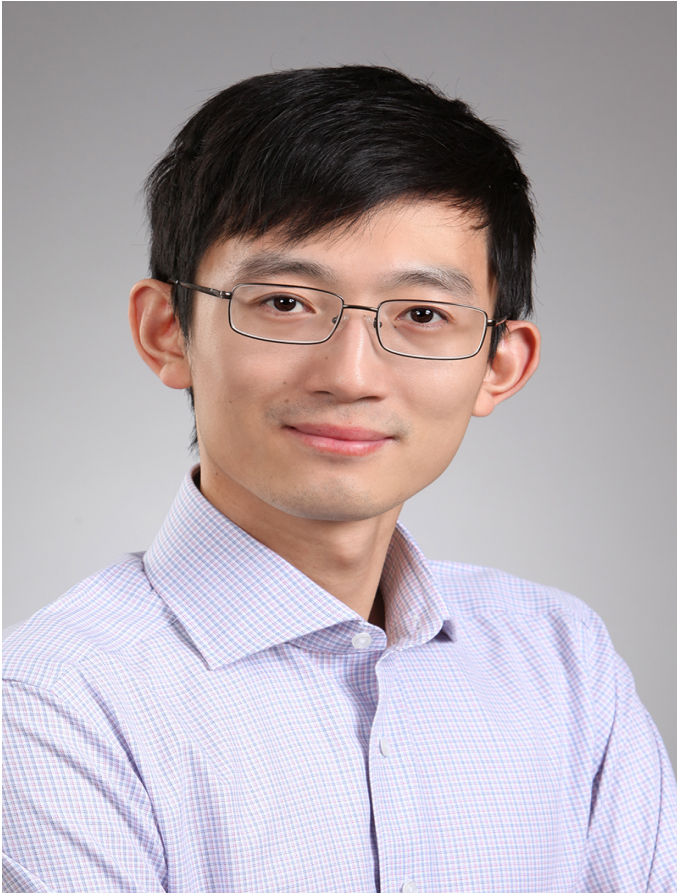
[“The University of Maryland’s Chinese-English Neural Machine Translation Systems at WMT18”](#). WMT 2018.

Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio.

[“Neural Machine Translation by Jointly Learning to Align and Translate”](#). ICLR 2015.

# MATHEMATICAL FRONTIERS

## Machine Learning for Text



**Tengyu Ma,  
Stanford University**



**Marine Carpuat,  
University of Maryland**



**Mark Green,  
UCLA (moderator)**

# MATHEMATICAL FRONTIERS

## 2019 Monthly Webinar Series, 2-3pm ET

**February 12:** *Machine Learning for Materials Science\**

**March 12:** *Mathematics of Privacy\**

**April 9:** *Mathematics of Gravitational Waves\**

**May 14:** *Algebraic Geometry\**

**June 11:** *Mathematics of Transportation\**

**July 9:** *Cryptography & Cybersecurity\**

**August 13:** *Machine Learning in Medicine\**

**September 10:** *Logic and Foundations\**

**October 8:** *Mathematics of Quantum Physics\**

**November 12:** *Quantum Encryption\**

**December 10:** *Machine Learning for Text*

*Made possible by support for BMSA from the  
**National Science Foundation**  
**Division of Mathematical Sciences**  
and the  
**Department of Energy**  
**Advanced Scientific Computing Research***

*\* Webinar posted*