

Neural Networks for Machine Translation

David Chiang

University of Southern California / University of Notre Dame

Joint work with

Victoria Fossum

Ashish Vaswani

Yingdong Zhao

Outline

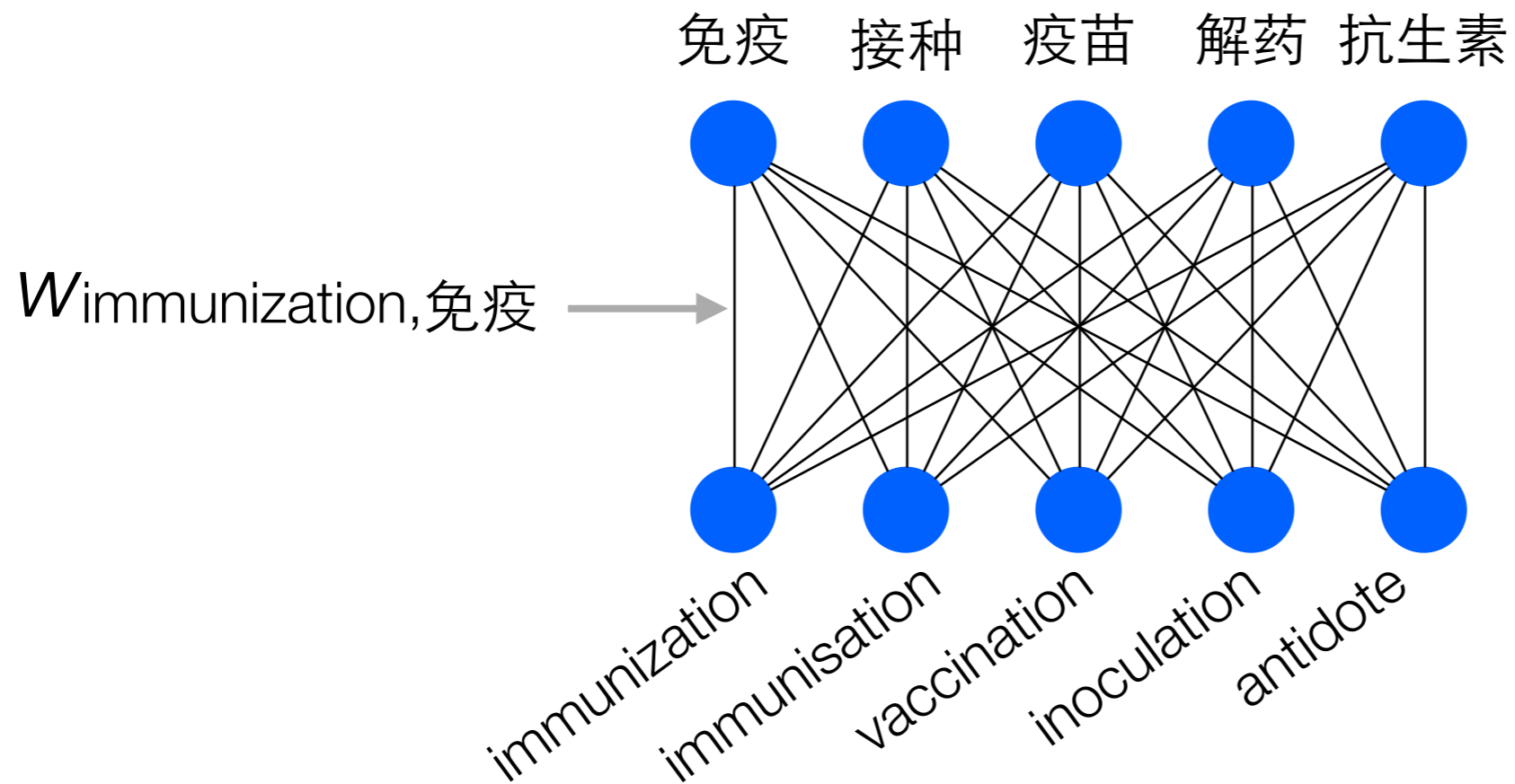
Introduction to neural networks

Fast training of neural networks

Neural networks for machine translation

Introduction to neural networks

Log-linear model



$$P(c|e) \propto \exp w_{e,c}$$

Word translation

English-to-Chinese

	免疫 miǎnyì	接种 jiēzhòng	疫苗 yìmiáo
immunization	0.04	0.1	0.04
immunisation	0.1	0.06	0.06
vaccination	0.06	0.2	0.2
inoculation	0	0.1	0.1

Features

$h(e) = 1$ iff...

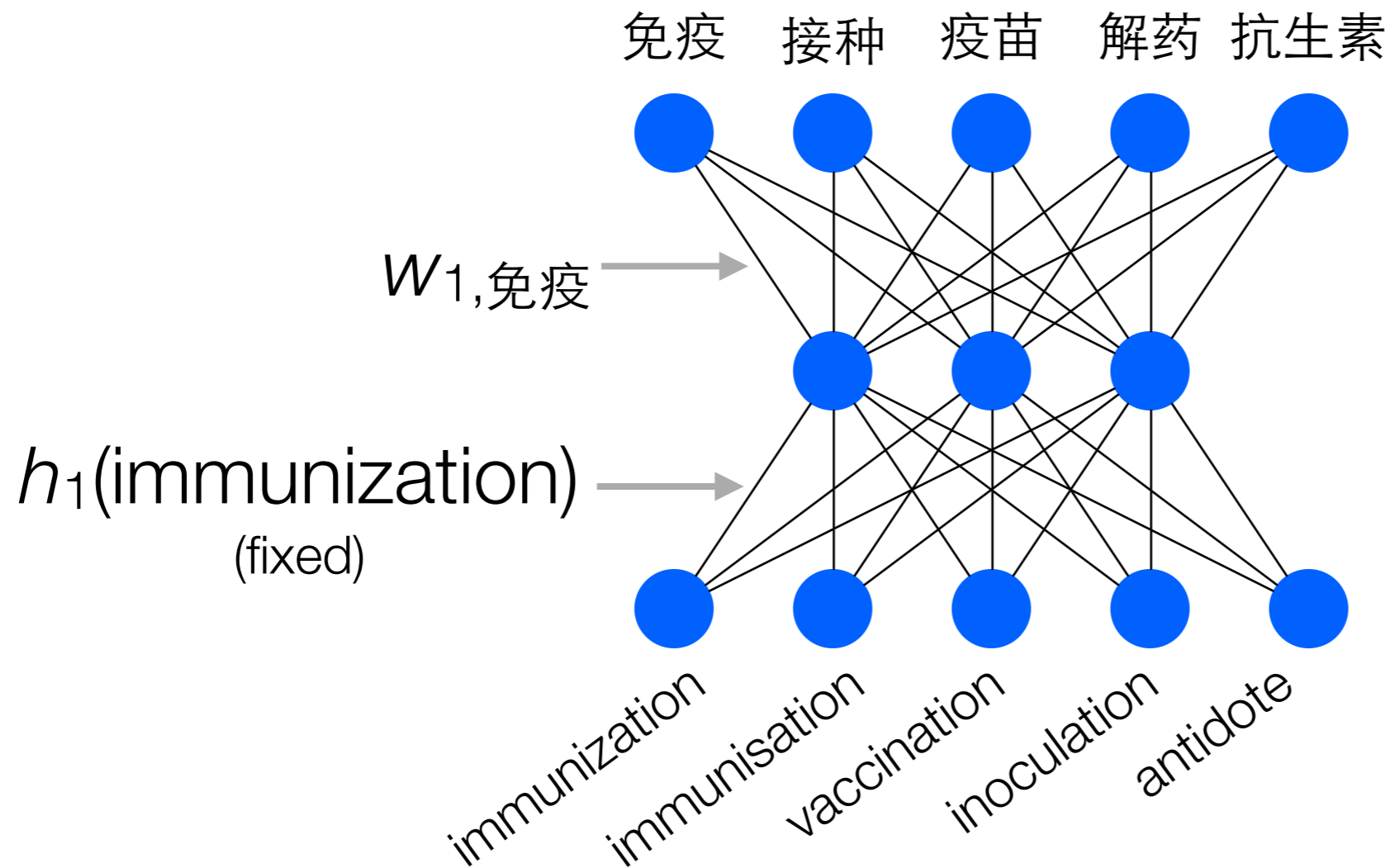
$e = \text{immunization}$

$\text{pos}(e) = \text{NN}$

$e[:5] = \text{immun-}$

$|e| < 5$

Log-linear model



$$P(c|e) \propto \exp \sum_i w_{i,c} h_i(e)$$

Conjunctions of features

$h(e) = 1$ iff...

$e = \text{immunization}$

$\text{pos}(e) = \text{NN}$

$e[:5] = \text{immun-}$

$|e| < 5$

$e = \text{immunization}$ and $\text{pos}(e) = \text{NN}$

$e = \text{immunization}$ and $e[-7:] = \text{-ization}$

$e = \text{immunization}$ and $|e| < 5$

$\text{pos}(e) = \text{NN}$ and $e[-7:] = \text{-ization}$

$\text{pos}(e) = \text{NN}$ and $|e| < 5$

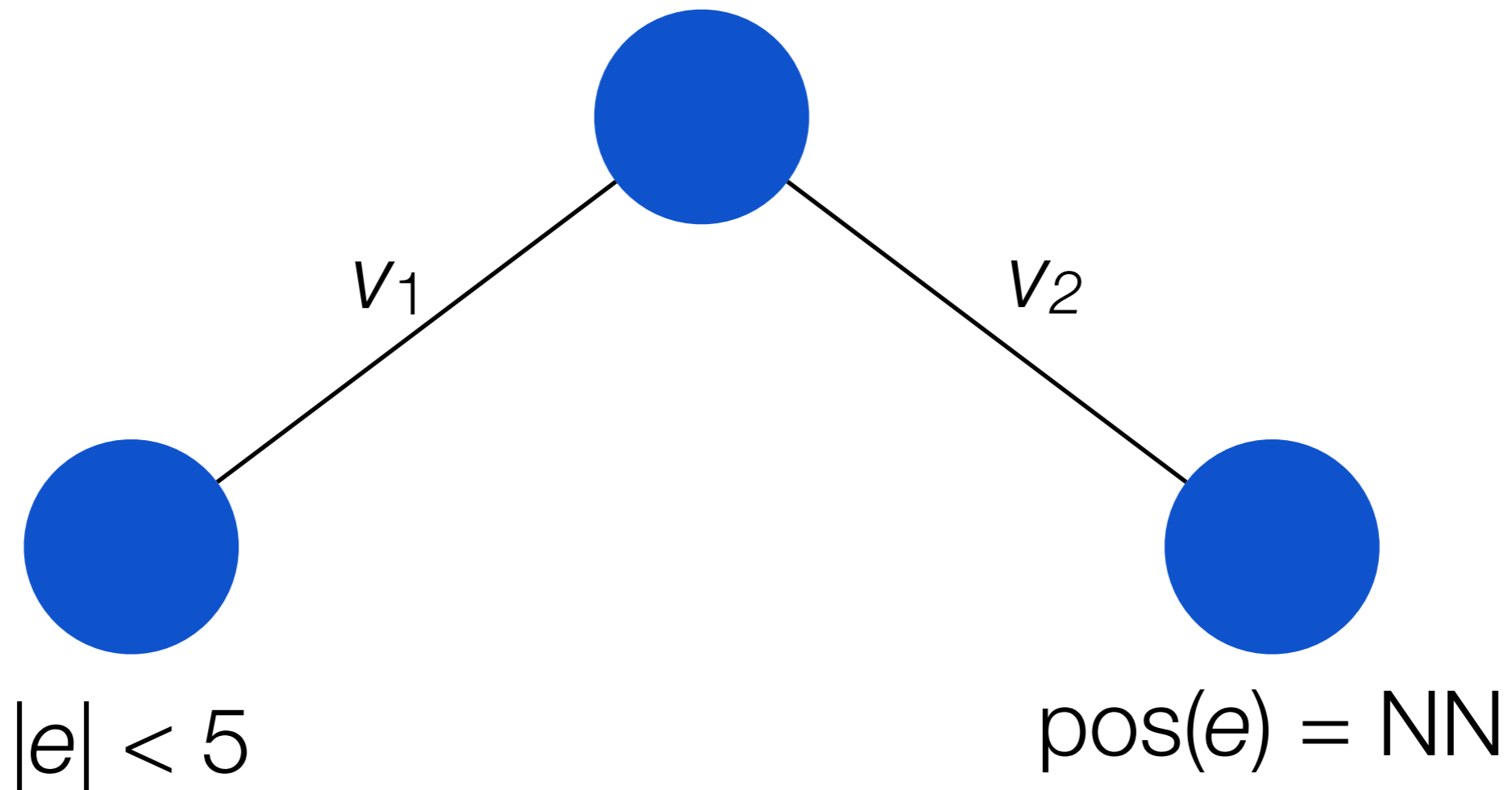
$e[-7:] = \text{-ization}$ and $|e| < 5$

“Programmable” features

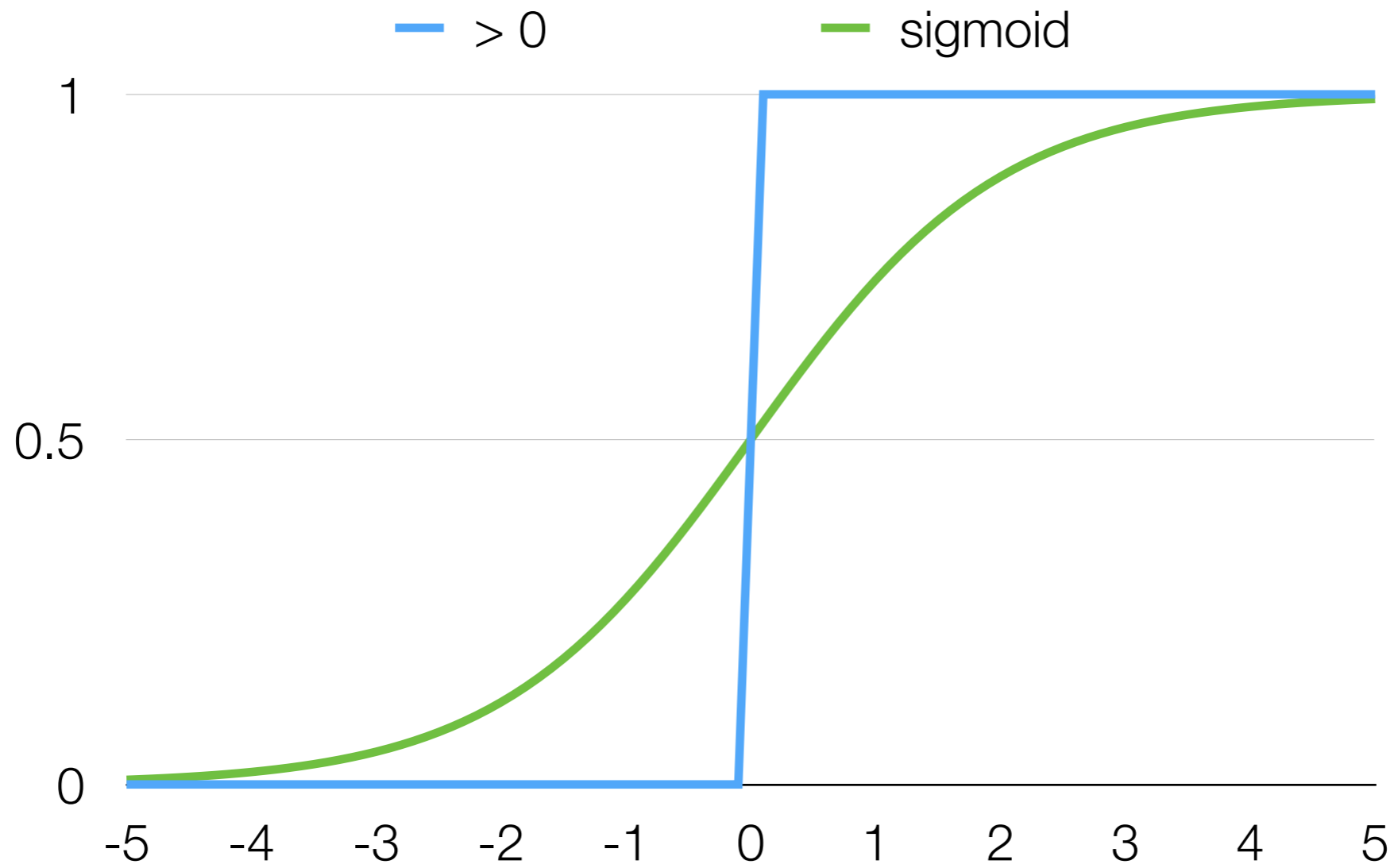
$\frac{1}{2} - (\text{pos}(e) = \text{NN}) > 0$	not $(\text{pos}(e) = \text{NN})$
$(e < 5) + (\text{pos}(e) = \text{NN}) - 1\frac{1}{2} > 0$	$ e < 5$ and $\text{pos}(e) = \text{NN}$
$(e < 5) + (\text{pos}(e) = \text{NN}) - \frac{1}{2} > 0$	$ e < 5$ or $\text{pos}(e) = \text{NN}$

“Programmable” features

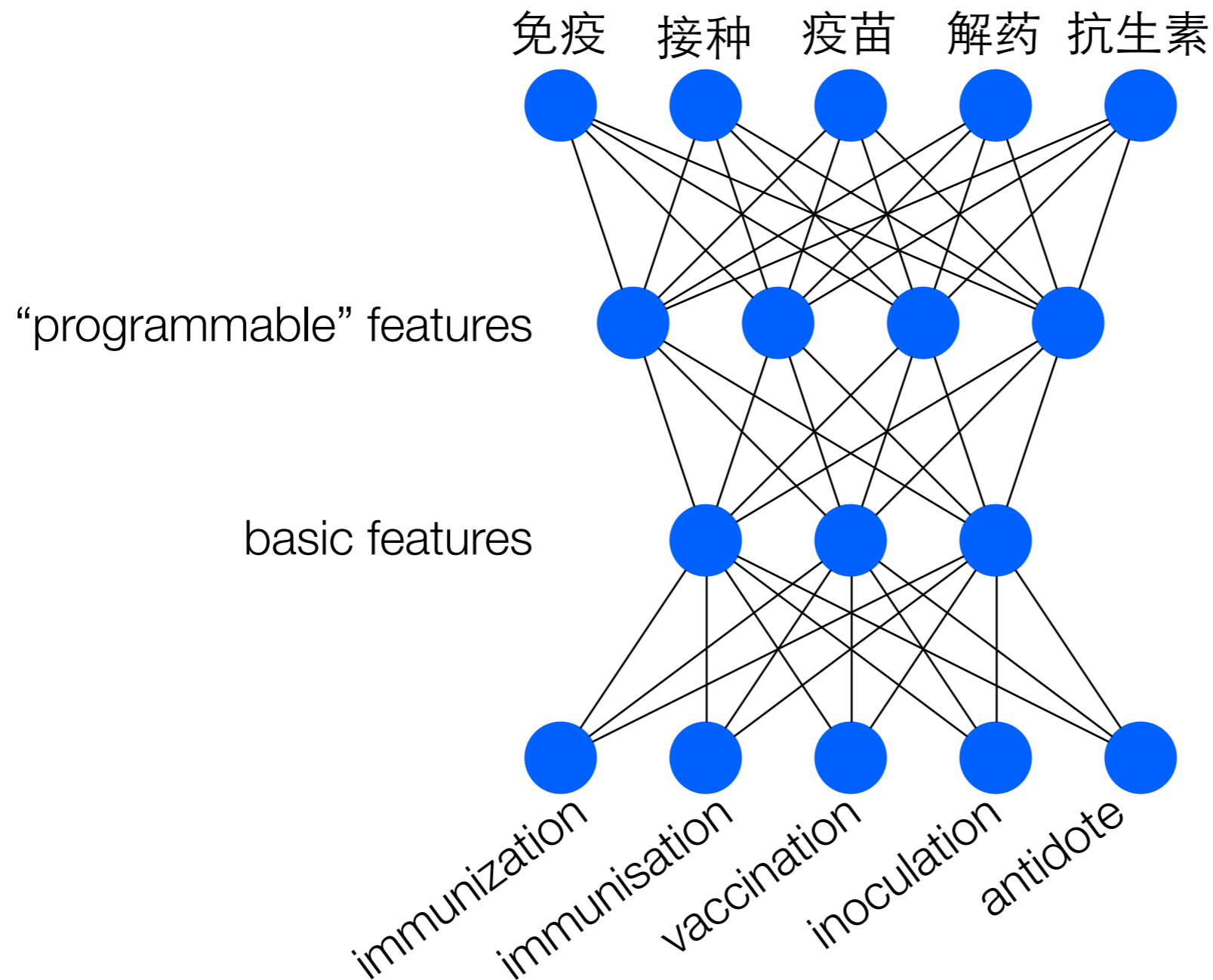
$$v_1(|e| < 5) + v_2(\text{pos}(e) = \text{NN}) + b > 0$$



Activation functions

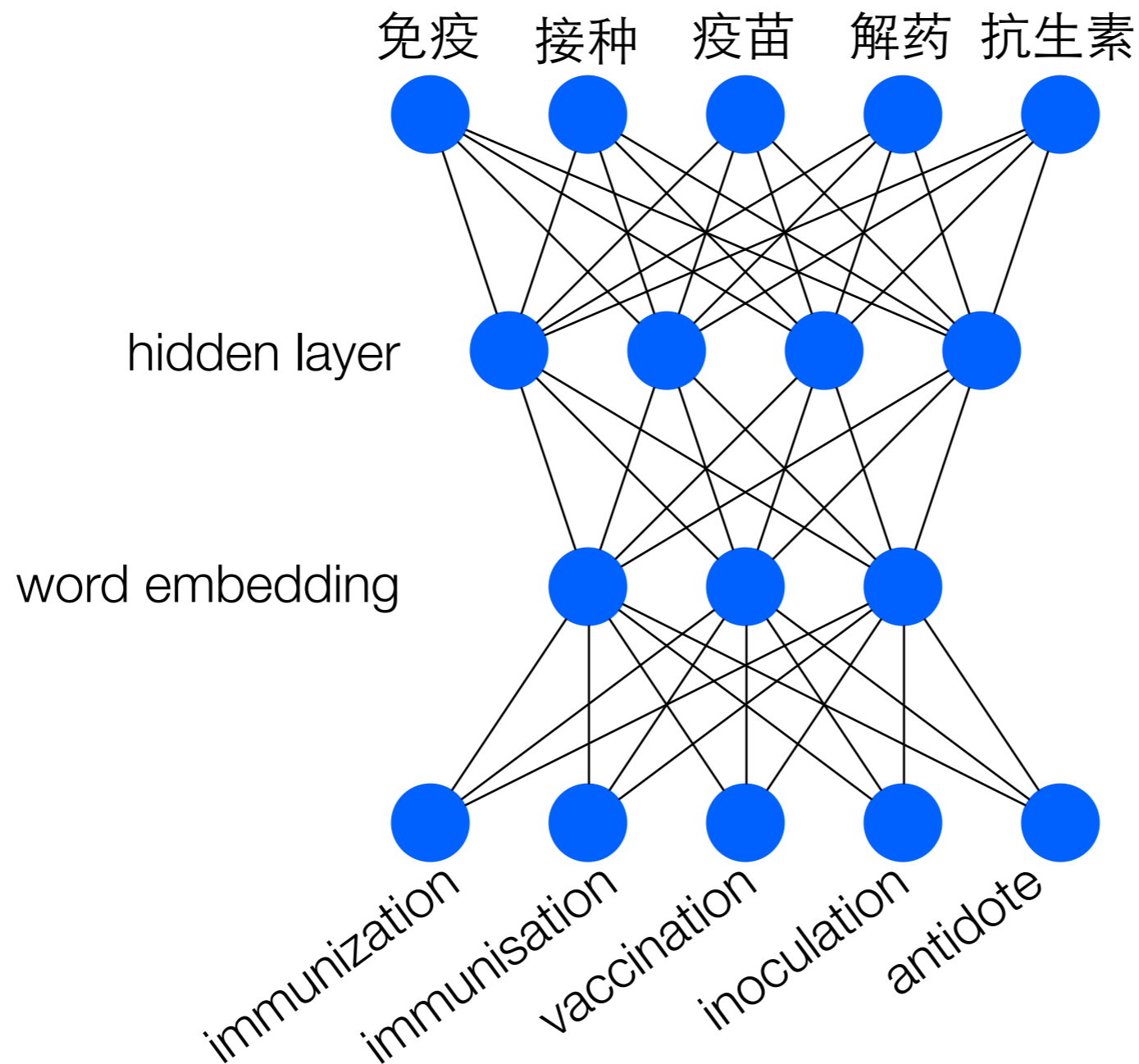


“Programmable” features as hidden layer



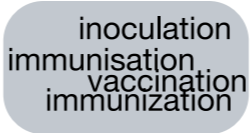
$$P(c|e) \propto \exp \sum_j w_{j,c} \text{sigmoid}(\sum_i v_{i,j} h_i(e) + b)$$

Feedforward neural network



$$P(c|e) \propto \exp \sum_j w_{j,c} \text{sigmoid}(\sum_i v_{i,j} u_{e,i} + b)$$

accommodation, orientation experimentation, cognition, vacation, affliction, desertion, delamination, extortion, immolation, decapitation, defection, assassination, molestation, mutilation, dislocation, disorientation, perspiration, laceration, indigestion, masturbation, fibrillation, menstruation, infarction, infliction, strangulation, amputation, castration, conviction, incarceration, abduction, execution, extradition, deportation, eviction, deportation, repatriation, relocation, dissection, urination, constipation, malformation, intoxication, suffocation, asphyxiation, electrocution, probation, investigation, interrogaton, prosecution, litigation, hospitalization, hospitalisation, degradation, inflammation, infection, dehydration, exhaustion, malnutrition, starvation, mutation, radiation, degeneration, dysfunctions, retardation, dysfunctions, recuperation, detoxification, pollination, adulteration, irradiation, incineration, oxidation, fermentation, contamination, malfunction, siltation, infestation, sedimentation, inundation, implosion, devastation, prostitution, abortion, contraception, sterilization, tuition, munitioin, vibration, congestion, insulation, precipitation, evaporation, condensation, pressurization, lubrication, suction, illumination, gasification, desalination, ammunitioin, ignition, ventilation, refrigeration, navigation, irrigation, sanitation, nutrition, education, migration, emigration, naturalization, immigration, arbitration



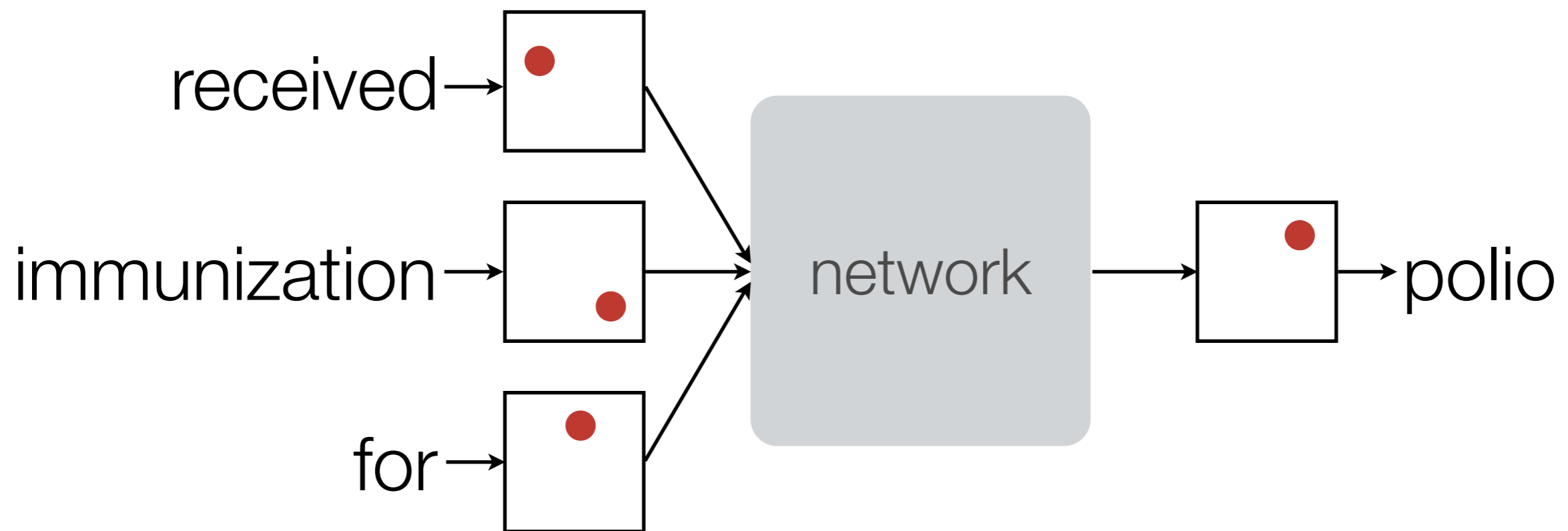
Fast training of neural networks

Ashish Vaswani, Yinggong Zhao, Victoria Fossum, and David Chiang. Decoding with large-scale neural language models improves translation. In *Proc. EMNLP 2013*.

Neural network language model

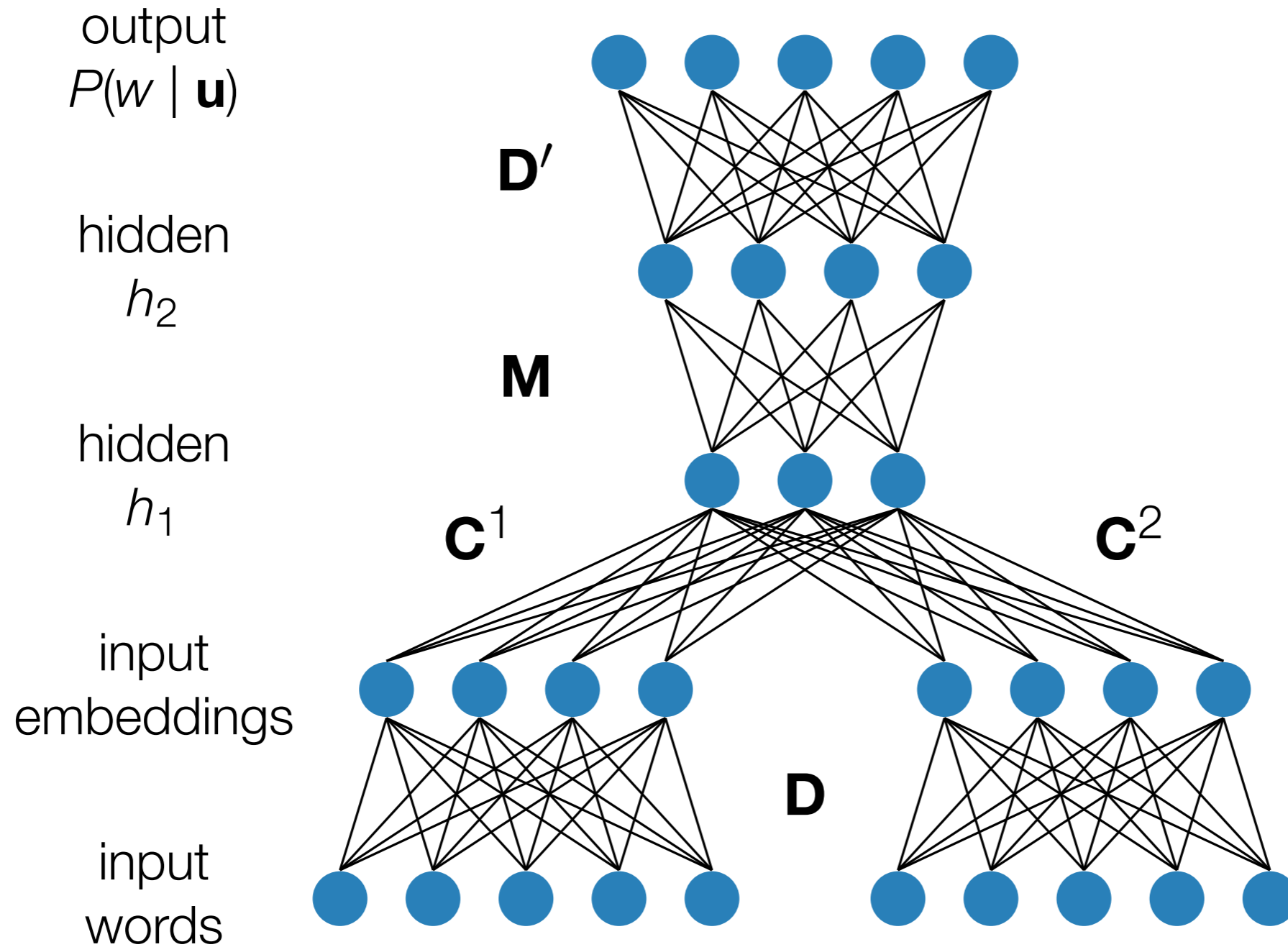
Bengio, 2003

Predict next English word w given previous English words \mathbf{u}



Neural network language model

Bengio, 2003



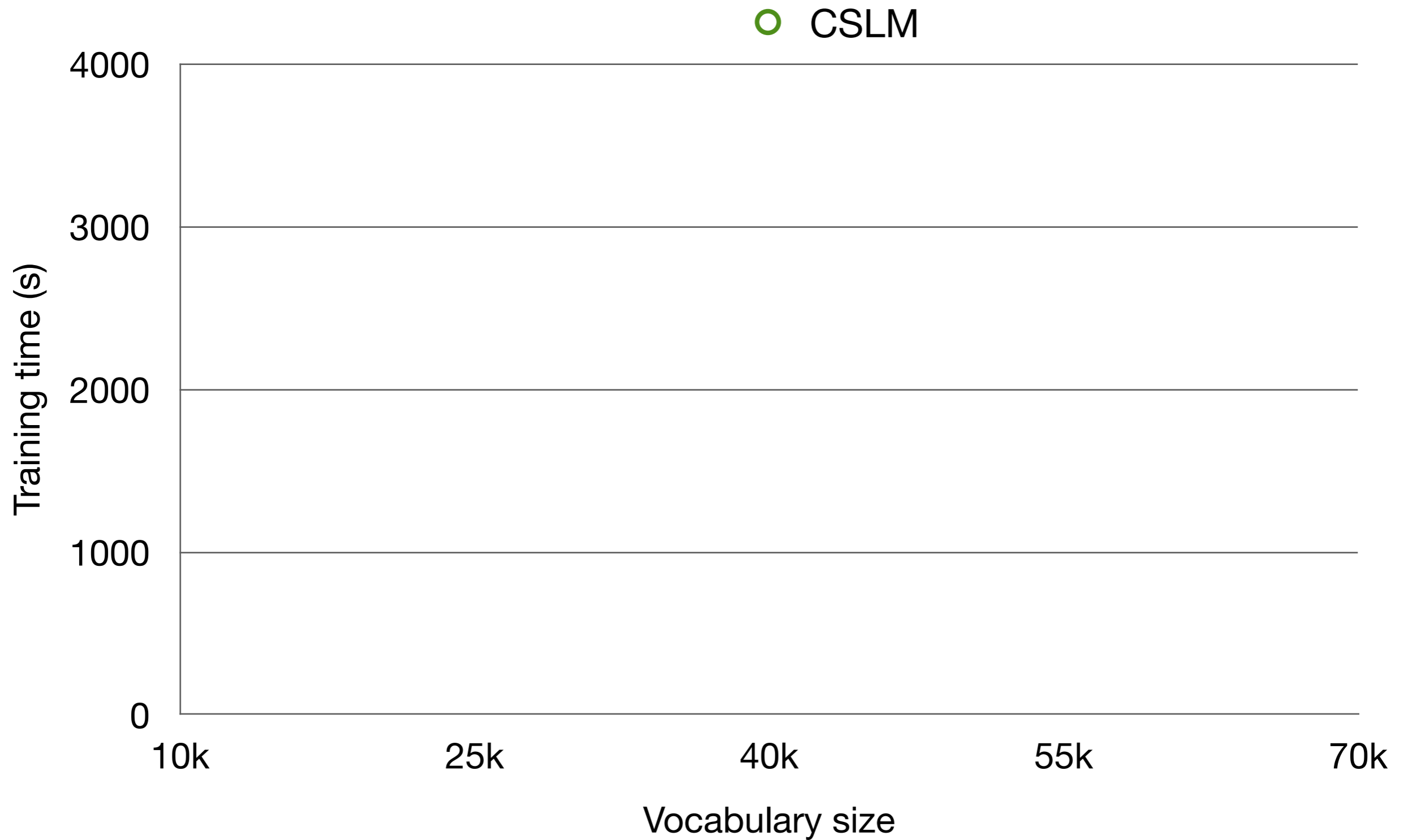
Neural network language model

Bengio, 2003

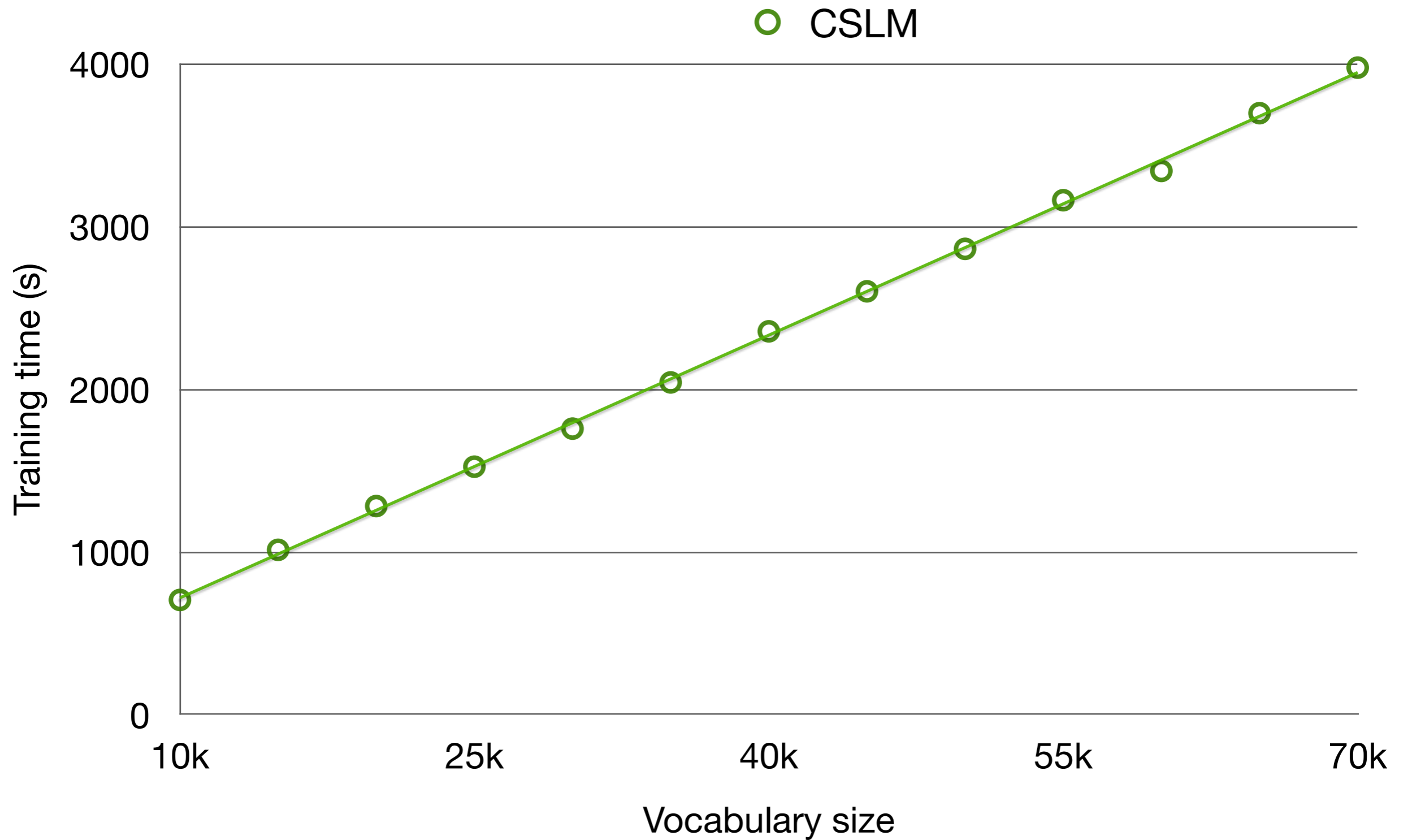
Problem: normalization factor Z makes both training and decoding slow

$$P(w | \mathbf{u}) = \frac{1}{Z(\mathbf{u})} \exp \dots$$
$$Z(\mathbf{u}) = \sum \exp \dots$$

Speed vs. vocabulary size

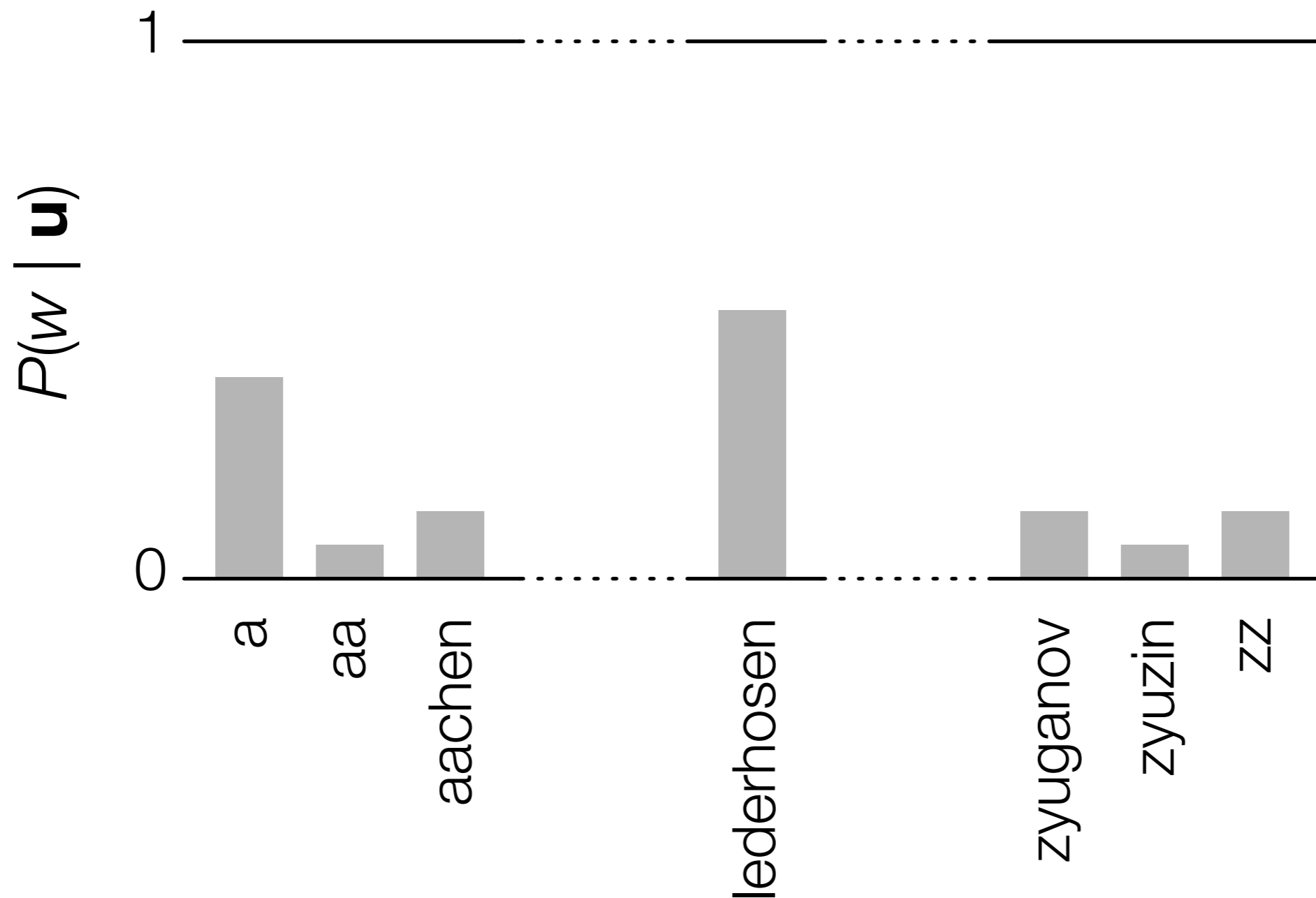


Speed vs. vocabulary size



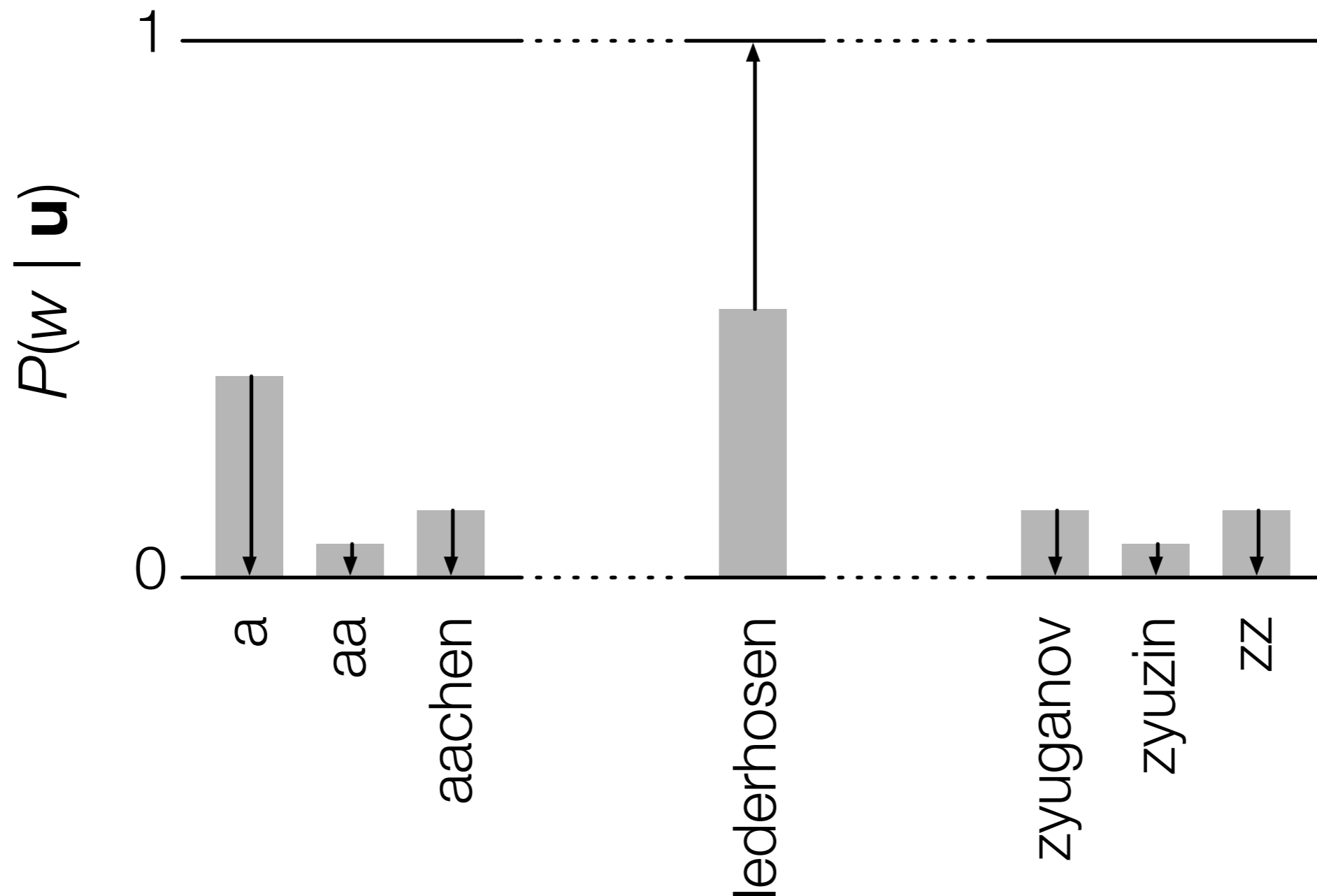
Maximum likelihood estimation

For each training example (\mathbf{u}, w) :



Maximum likelihood estimation

For each training example (\mathbf{u}, w) :



Noise contrastive estimation

Gutmann and Hyvarinen, 2010

For each training example (\mathbf{u} , w):

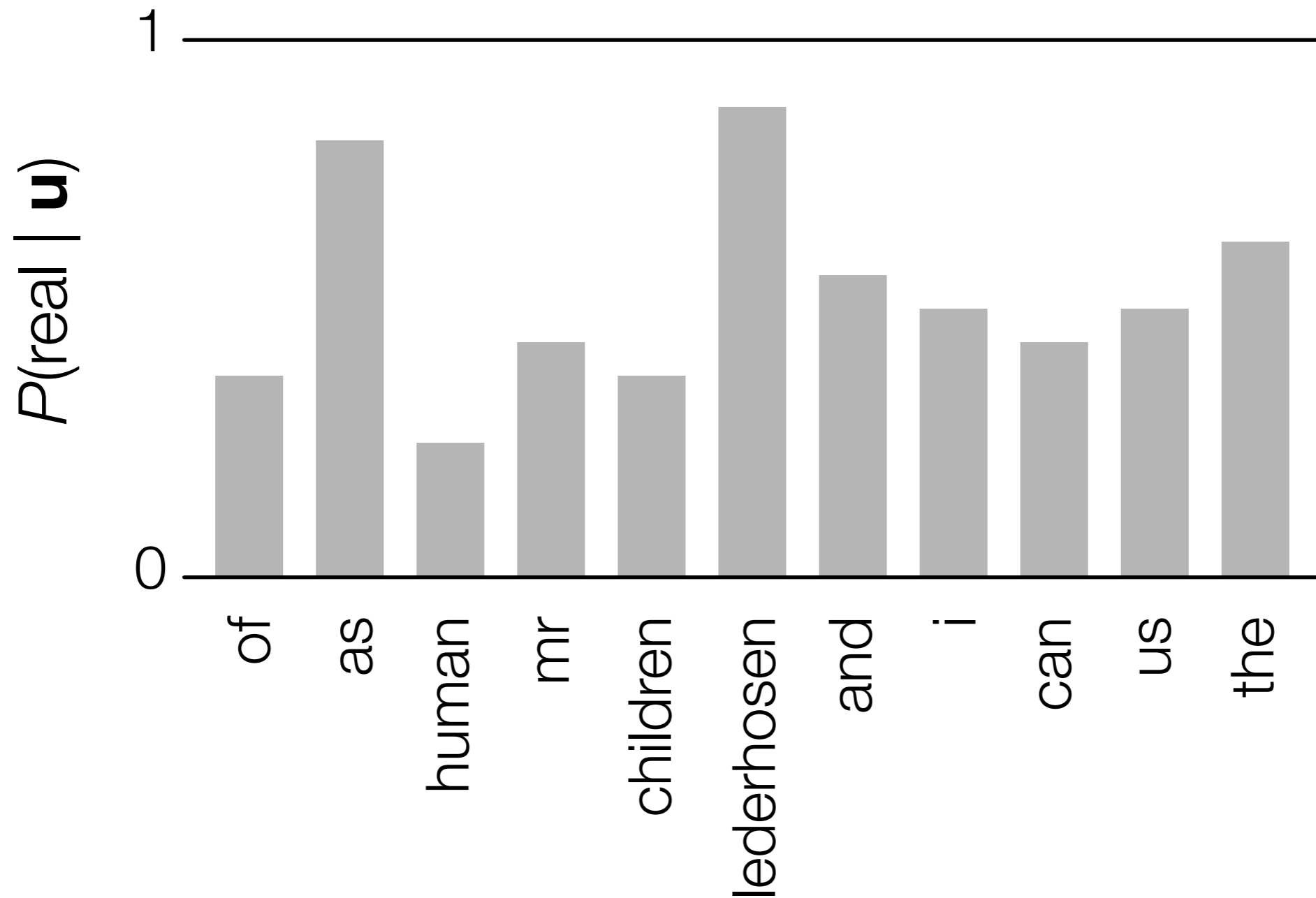
generate k noise samples

train model to classify real examples and noise samples

Noise contrastive estimation

Gutmann and Hyvarinen, 2010

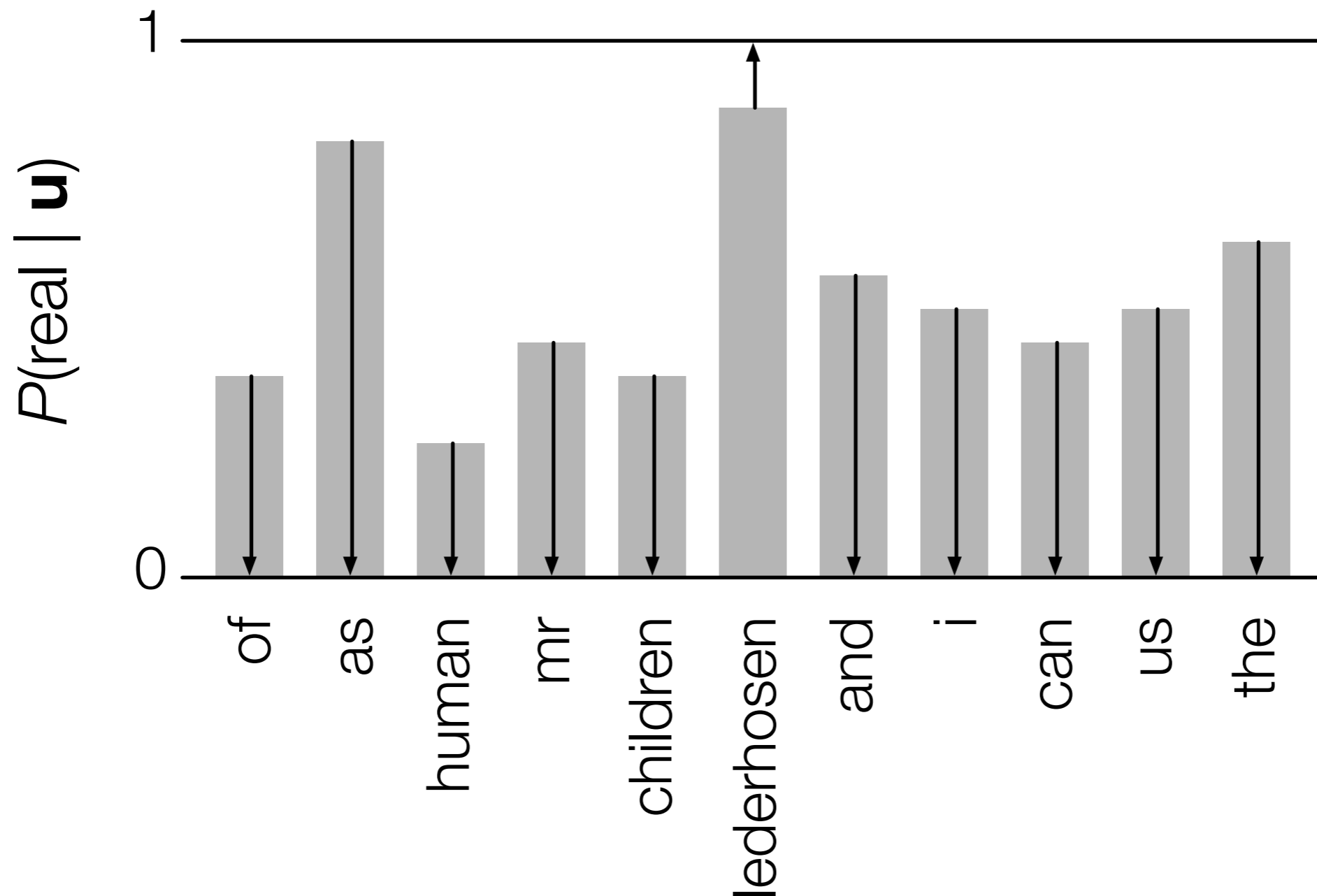
For each training example (\mathbf{u}, w) :



Noise contrastive estimation

Gutmann and Hyvarinen, 2010

For each training example (\mathbf{u}, w) :



Noise contrastive estimation

Gutmann and Hyvarinen, 2010

Bonus: this estimates Z also

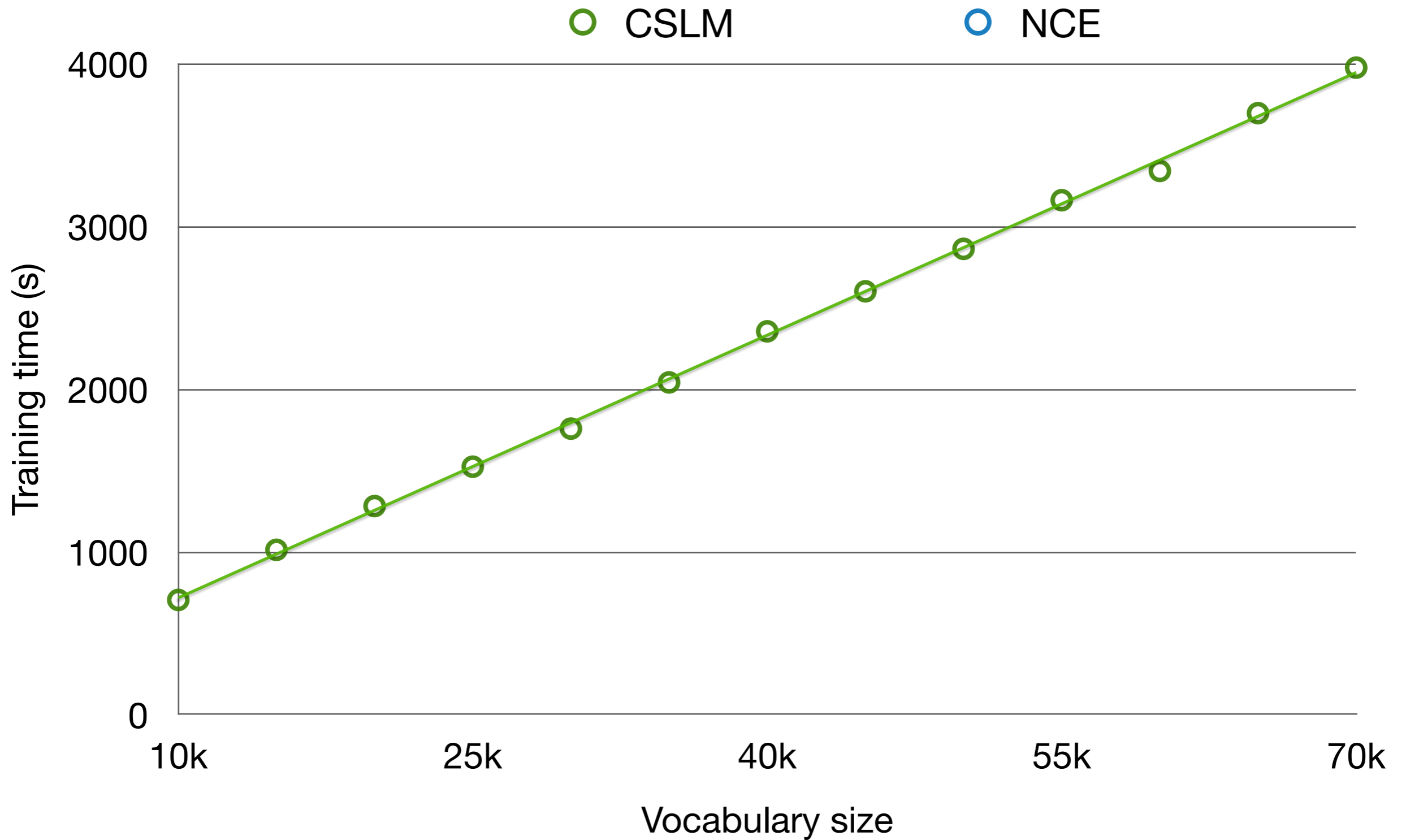
$$P(w | \mathbf{u}) = \frac{1}{Z(\mathbf{u})} \exp \dots$$

$$Z(\mathbf{u}) = \sum \exp \dots$$

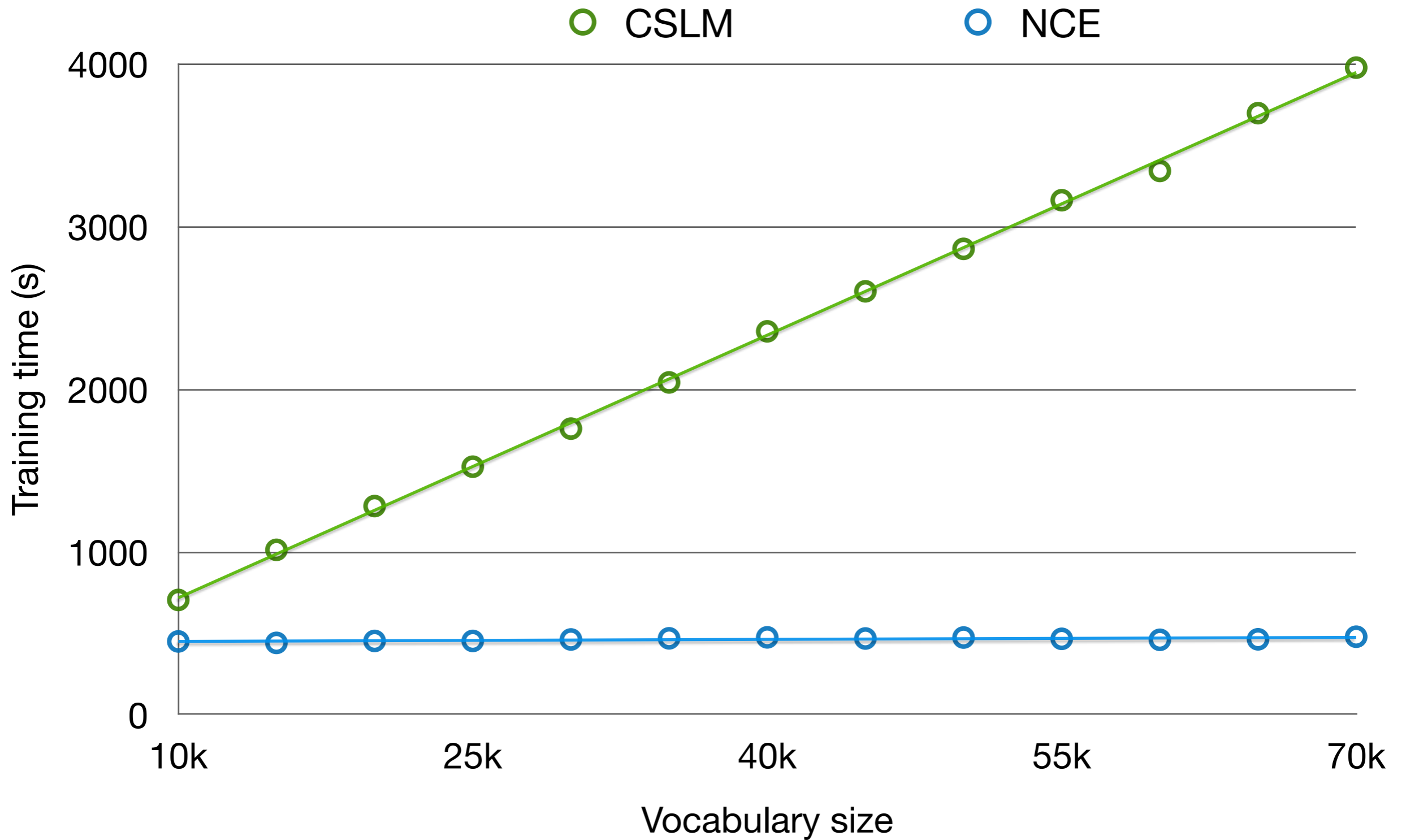
Mnih and Teh: just freeze $Z = 1$

⇒ During decoding, we can skip expensive calculation of Z

Speed vs. vocabulary size



Speed vs. vocabulary size

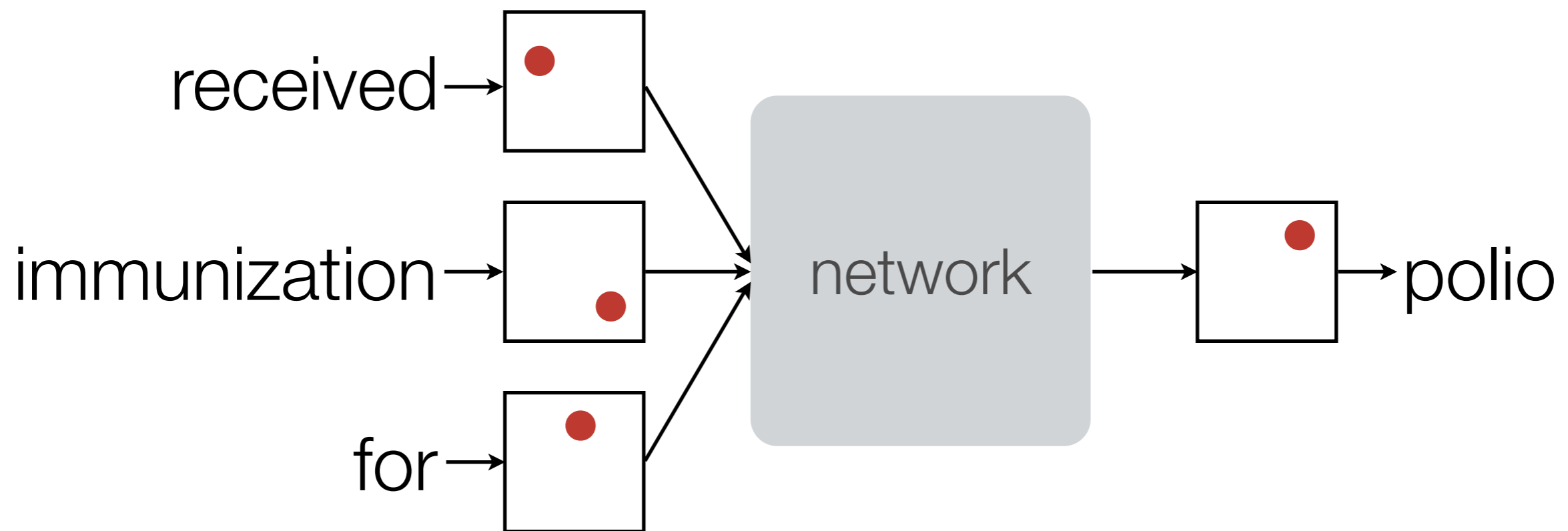


Neural networks for machine translation

Neural network language model

Bengio, 2003

Predict next English word w given previous English words \mathbf{u}



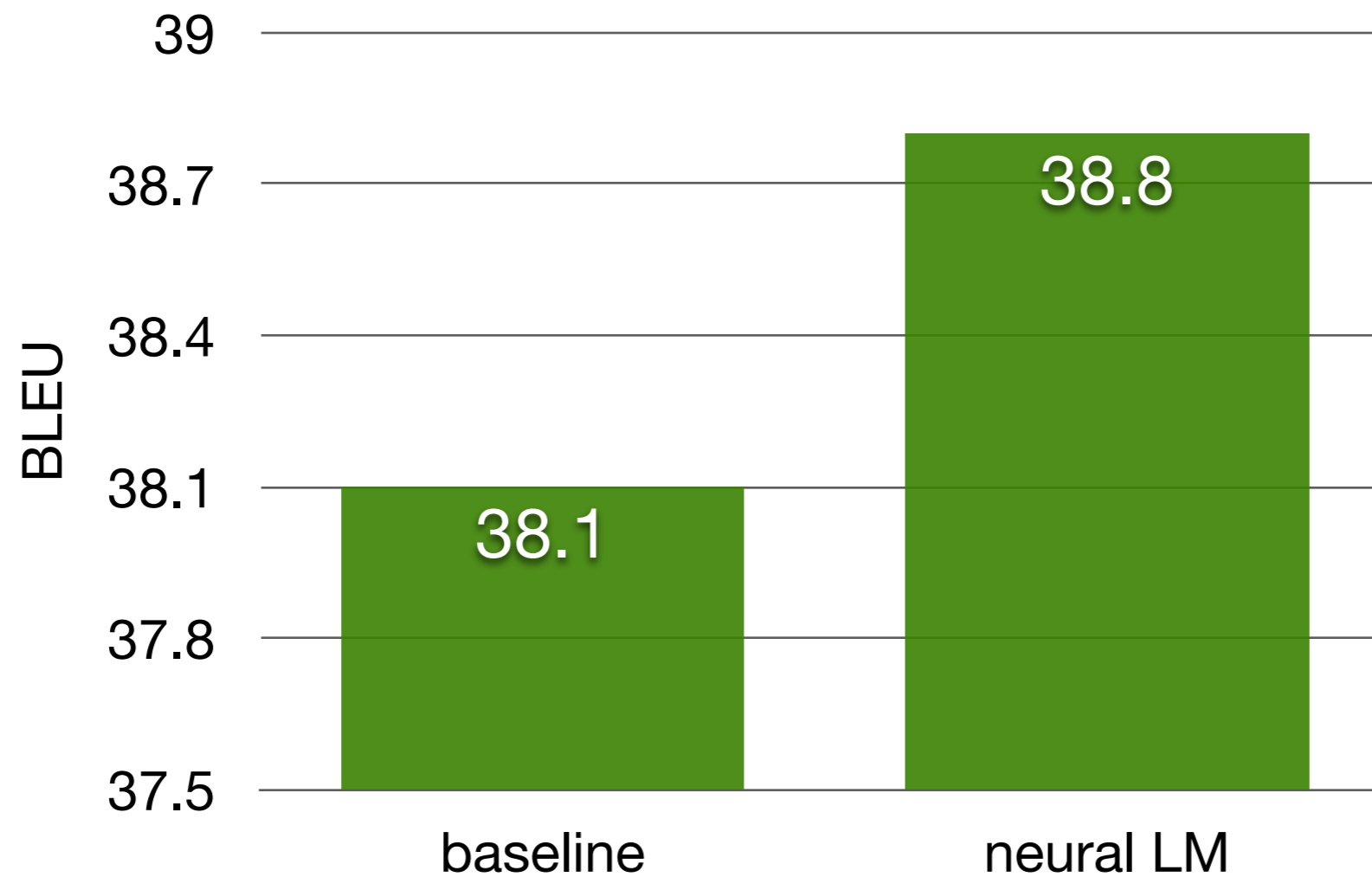
Translation results

Chinese-English BLEU

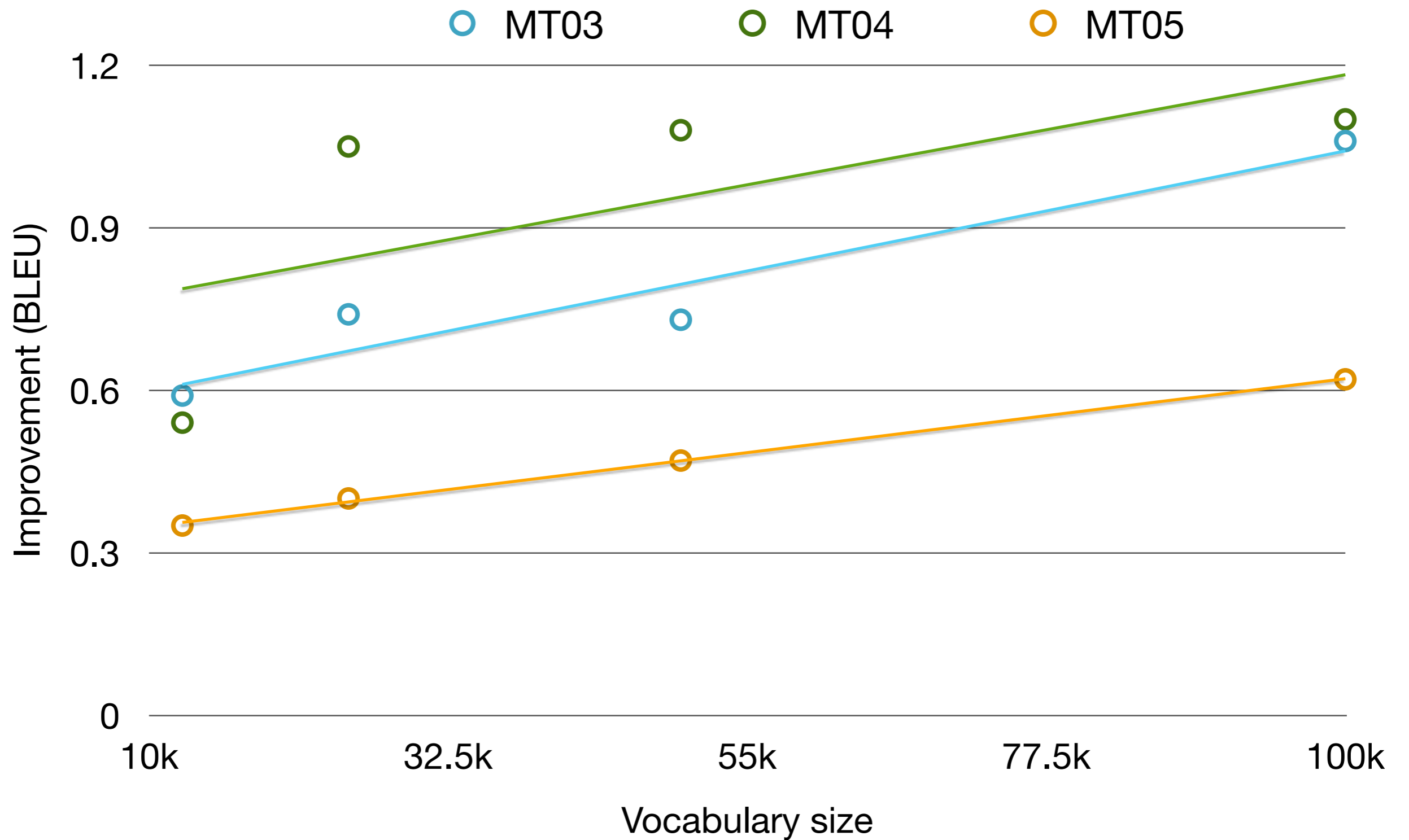
	2004	2005	2006
baseline	38.4	37.7	34.3
neural LM <i>n</i> -best reranking	38.6	37.8	34.7
neural LM during decoding	39.5	38.8	34.9

Translation results

BOLT Chinese-English forum



BLEU vs. Vocabulary size



Neural network translation models

Devlin et al., 2014

Predict next English word given Chinese words and previous English words

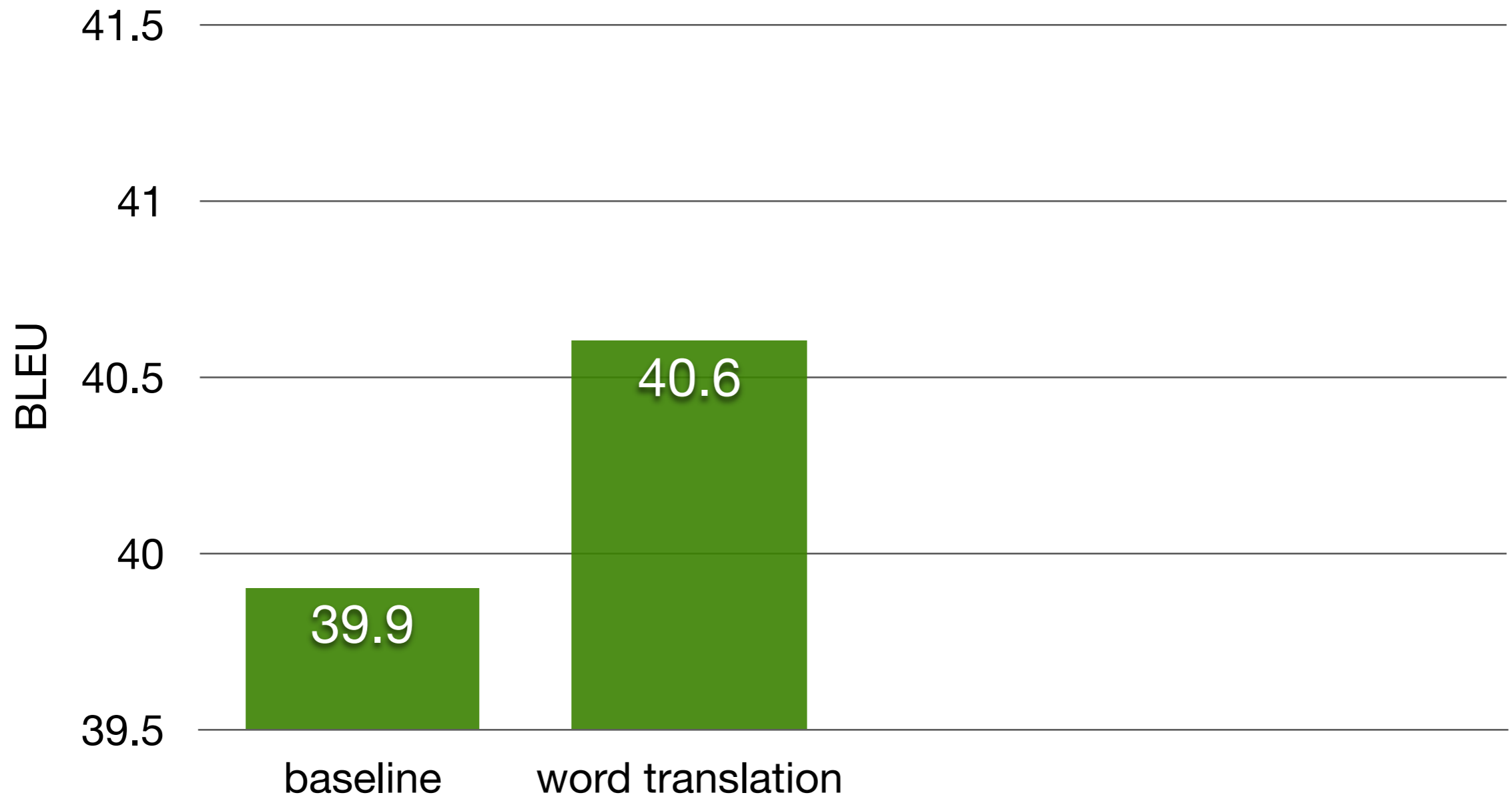
...that have diplomatic relations with North...

$P(\text{relations} \mid \text{帮教}, \dots)$

澳洲是与北韩有帮教的少数国家之一

Translation results

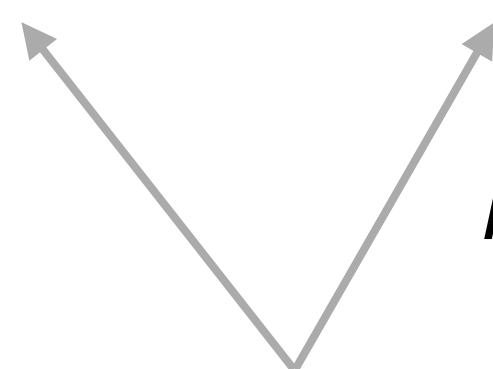
BOLT Arabic-English forum



Neural network translation models

Fertility: Predict number of English words
given Chinese words

...that have diplomatic relations with North...



$P(2 \mid \text{帮教}, \dots)$

澳洲是与北韩有帮教的少数国家之一

Neural network translation models

Distortion: Predict reordering given Chinese words and previous English words

...that have diplomatic relations with North...

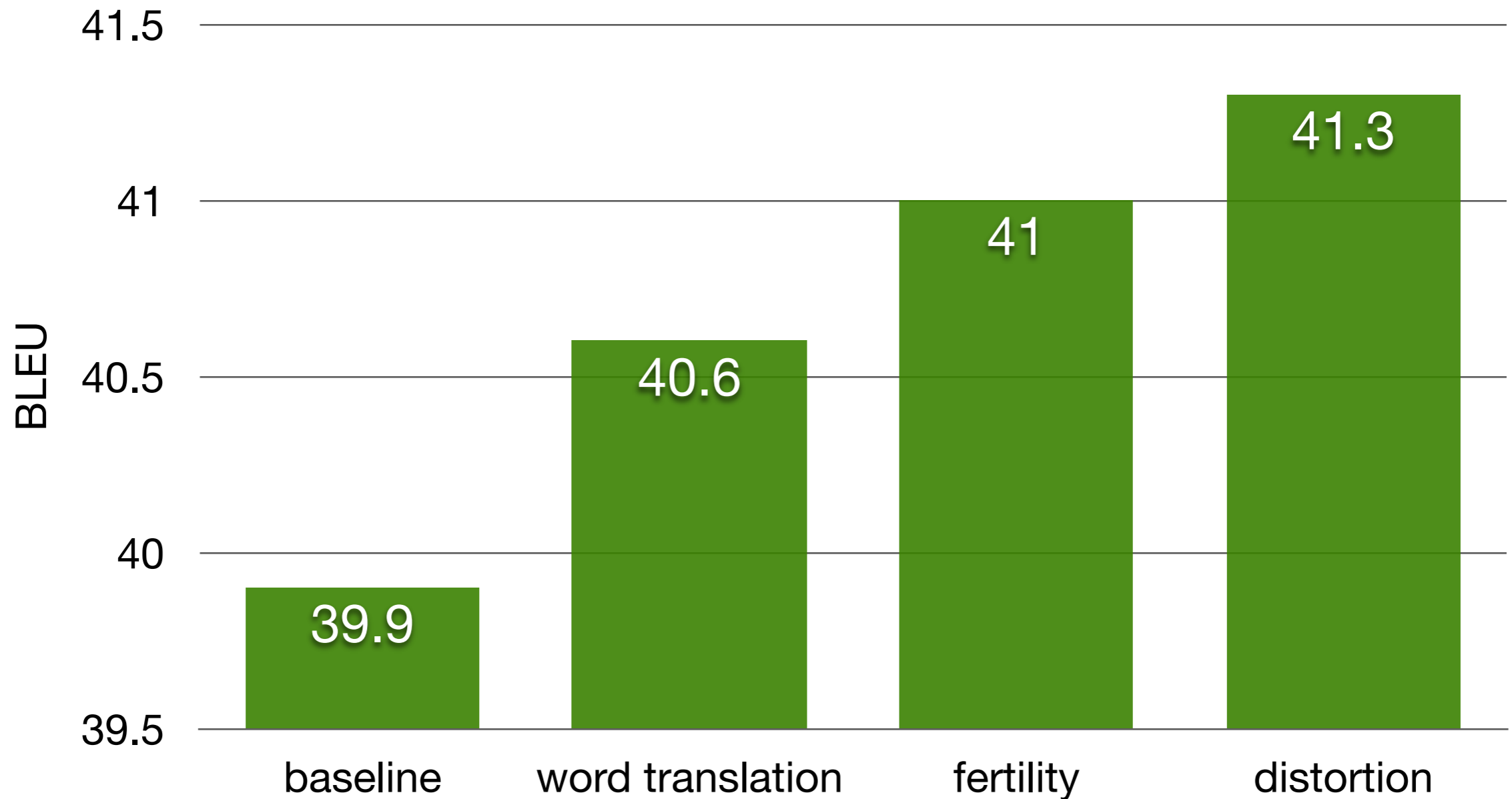
$P(-3 \mid \text{帮教}, \dots)$

澳洲是与北韩有帮教的少数国家之一

-3 -2 -1

Translation results

BOLT Arabic-English forum



Conclusions

- Neural networks can be thought of as automatically inventing features, even for very large contexts
- A major obstacle (large output vocabulary) is removed
- Deep networks, pretraining, etc. are not needed?
- Large improvements are possible