# Natural Language Processing with Small Feed-Forward Networks

Jan A. Botha    **Emily Pitler**    Ji Ma    Anton Bakalov

Alex Salcianu    David Weiss    Ryan McDonald    Slav Petrov

emnlp 2017
Copenhagen

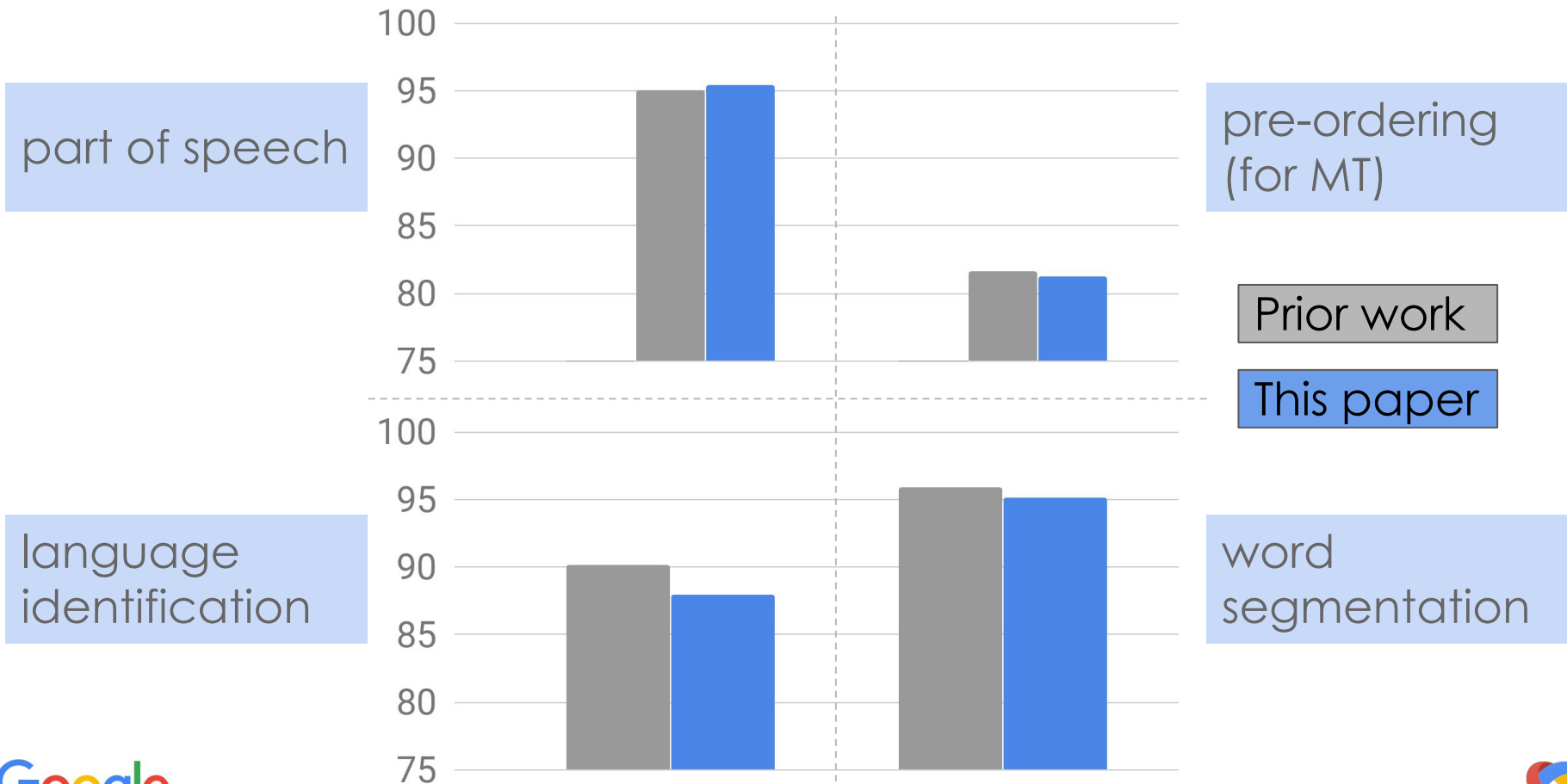# Goal: Small and Fast (on CPU/mobile)

Deep recurrent models: can have 100s of millions (or billions) of parameters

Recent work on smaller recurrent models: Kim and Rush, 2016; Sharp Models on Dull Hardware, Devlin, EMNLP 2017 10's to 100 tokens/second
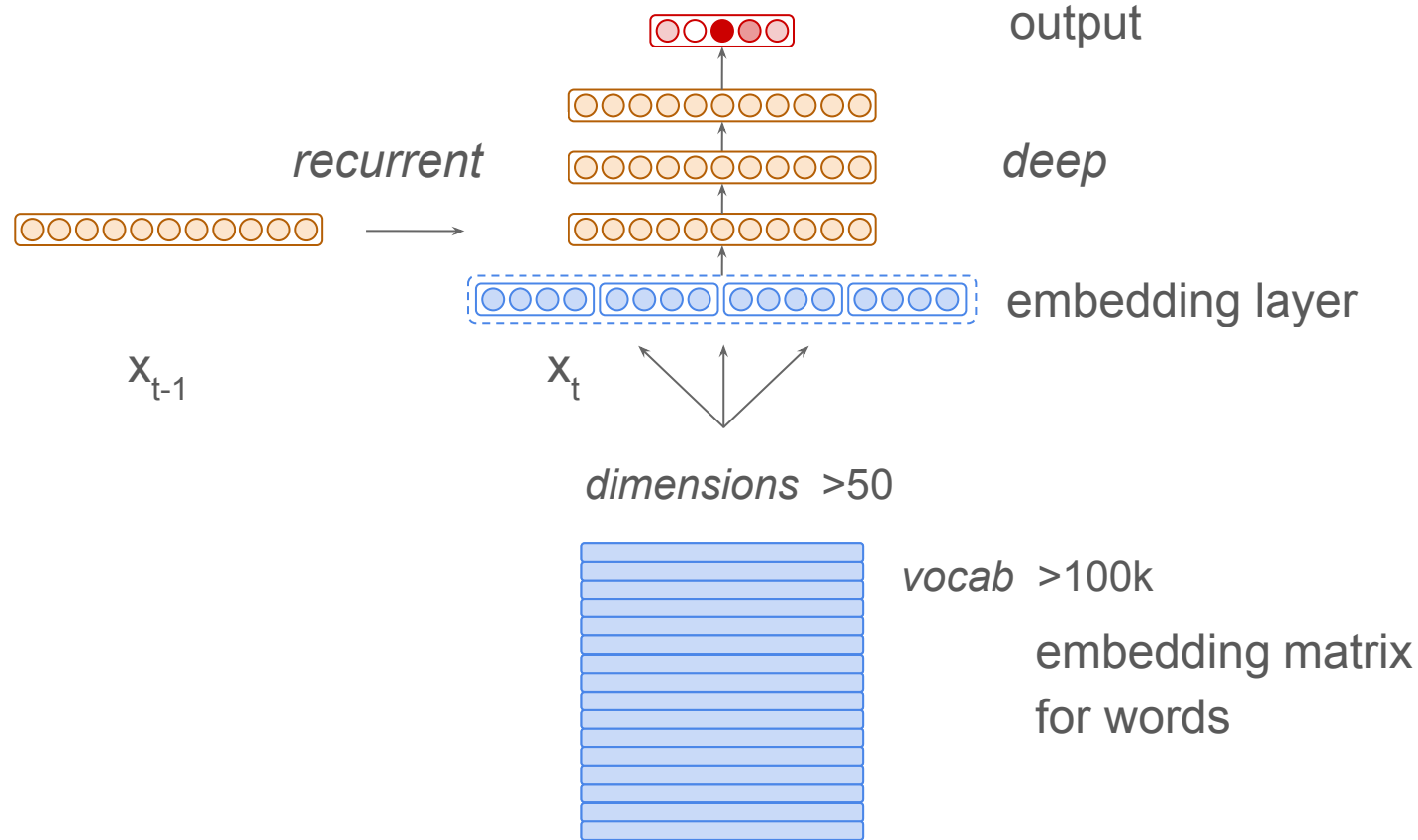
For variety of NLP tasks, can get order of magnitude speedup over LSTMs

| Memory | ≤ 2 MB |
| ---: | :--- |
| Speed | 7k – 46k tokens/second |
| Trained | in a few hours |

Google

# ...With Near State-of-the-Art Accuracies in 4 Tasks

part of speech

pre-ordering
(for MT)

Prior work

This paper

language
identification

word
segmentation

Google

# A Recurrent & Deep Model with Large Vocabulary

# A Recurrent & Deep Model with Large Vocabulary

output

recurrent

deep

**Slower Speed**

embedding layer

$x_{t-1}$

$x_t$

**Bigger Size**

dimensions >50

vocab >100k

embedding matrix for words

Google

# How Best to Allocate a Small Memory Budget?

output

*recurrent*

*deep*

embedding layer

$x_{t-1}$

$x_t$

~~embedding matrix~~
~~for words~~

Google

# What's Left? Small Feed-Forward Architecture

softmax
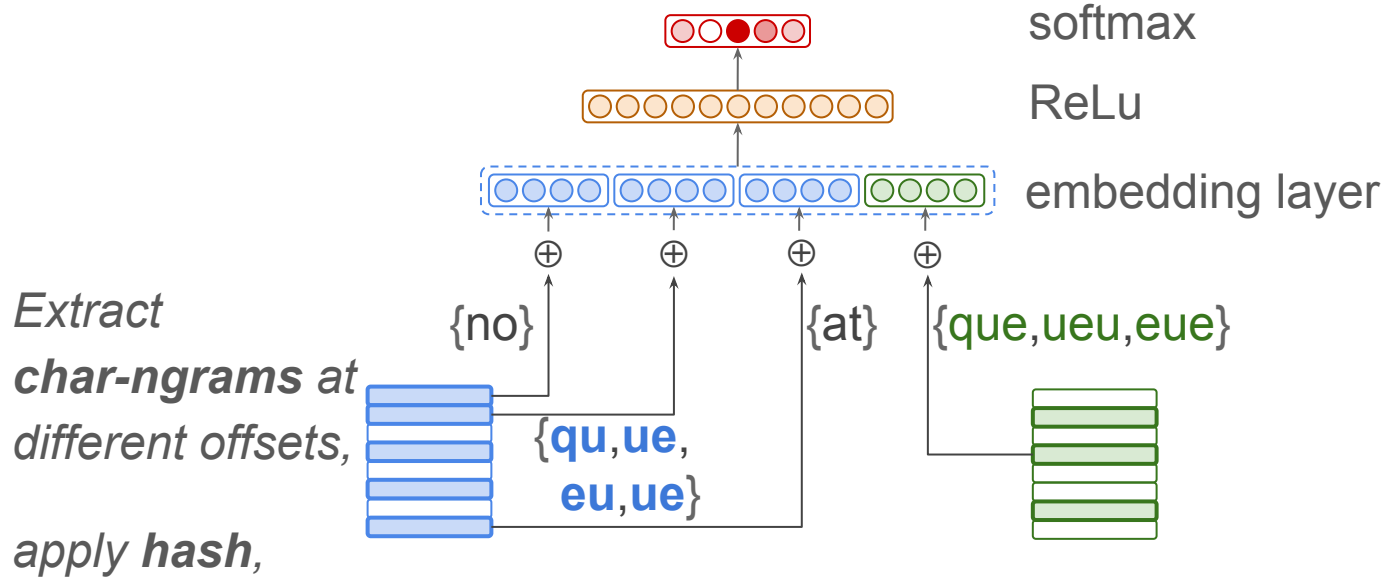
ReLu

embedding layer

Goals

⬆ Accuracy
─────────
⬇ Model Size

Word Clusters
Pipelines
Selected Features

Hashed Character n-grams
Quantization

Google

# Input Representation: Hashed Character n-grams

softmax

ReLu

embedding layer

⊕ ⊕ ⊕ ⊕

Extract **char-ngrams** at different offsets,

apply **hash**,

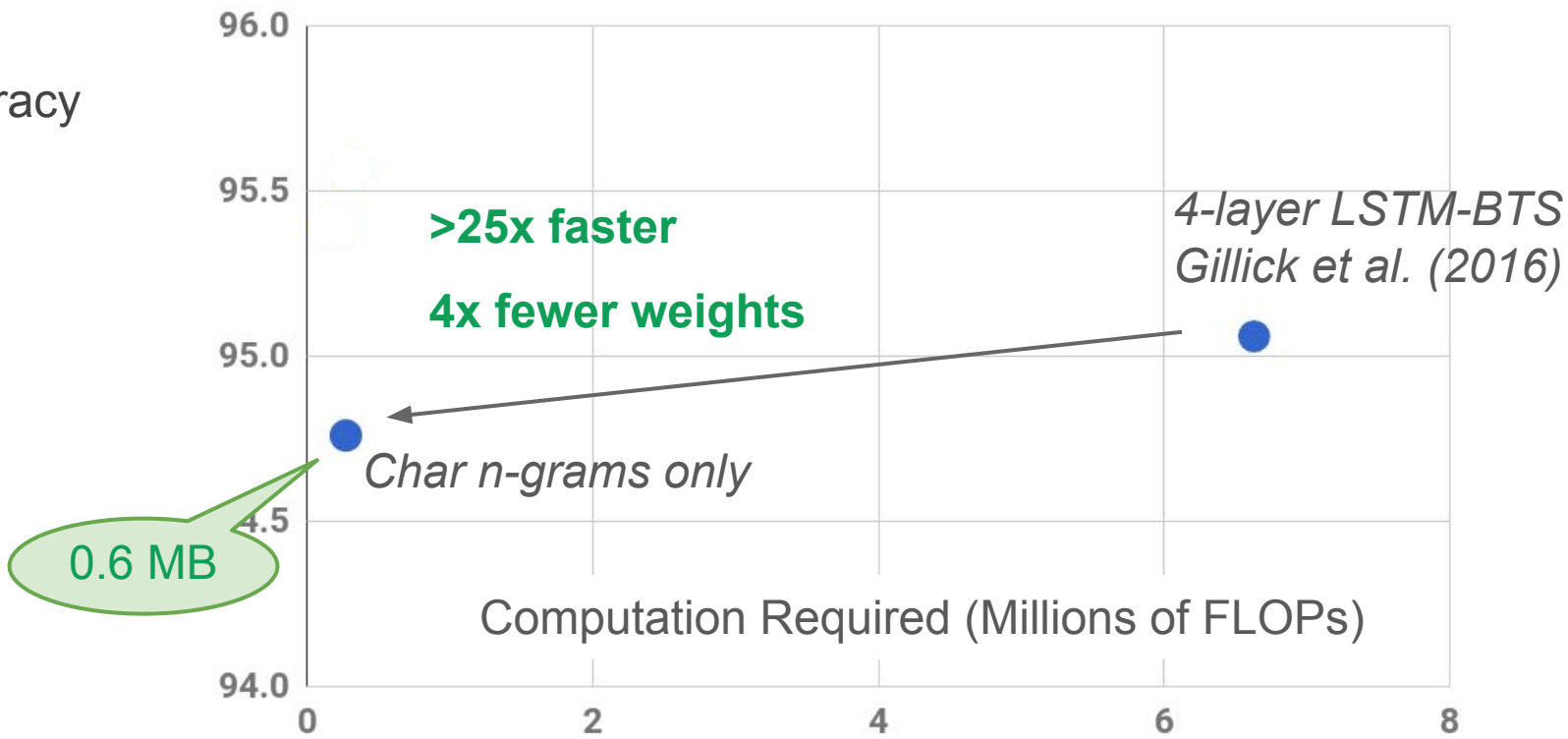aggregate looked up **embeddings**.

{no}

{qu,ue,
eu,ue}

{at}

{que,ueu,eue}

There was no **queue** at the ...

# Case Study 1: POS Tagging

# Vanilla Model: Less Resources, Little Less Accurate



**POS** Accuracy

96.0

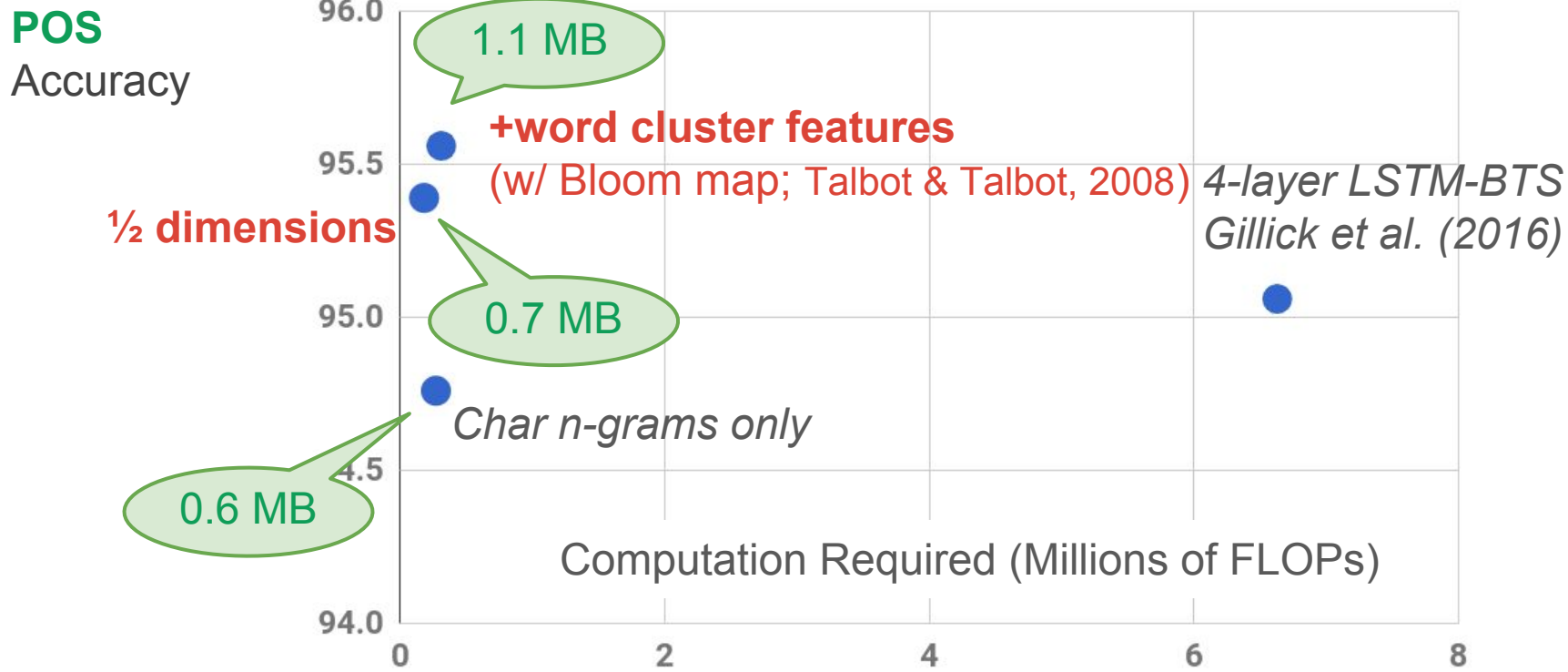95.5

95.0

4.5

94.0

>25x faster

4x fewer weights

*4-layer LSTM-BTS*
*Gillick et al. (2016)*

*Char n-grams only*

0.6 MB

Computation Required (Millions of FLOPs)

0    2    4    6    8

Faster ←——————————— Slower

Google

# Accuracy Boost Adding Resource-backed Features

**POS**
Accuracy



1.1 MB

**+word cluster features**
(w/ Bloom map; Talbot & Talbot, 2008)

*4-layer LSTM-BTS*
*Gillick et al. (2016)*

*Char n-grams only*

0.6 MB

Computation Required (Millions of FLOPs)

Faster ⟵                    Slower

# Allows Reducing Embedding Dimensions Further



POS Accuracy

**+word cluster features** (w/ Bloom map; Talbot & Talbot, 2008)

½ dimensions

1.1 MB

0.7 MB

0.6 MB

Char n-grams only

4-layer LSTM-BTS Gillick et al. (2016)

Computation Required (Millions of FLOPs)

Faster ← → Slower

Google

# Case Study 2: Preordering for MT

# Transition System for Structured Output

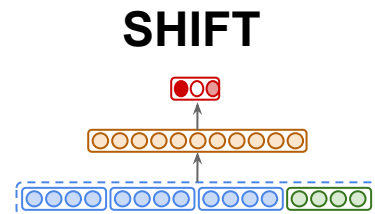English → Japanese word ordering: I ate pizza → I pizza ate

**Stack**
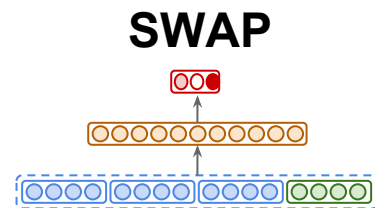
**Buffer**

**SWAP**

I ate pizza

I pizza

ate

**SHIFT**

Google

# POS Tags as Intermediate Representation

**Stack**

**Buffer**

**SHIFT**

I/PRP pizza/NN

ate/VBD

**PRP**

**NN**

**VBD**

I ate pizza
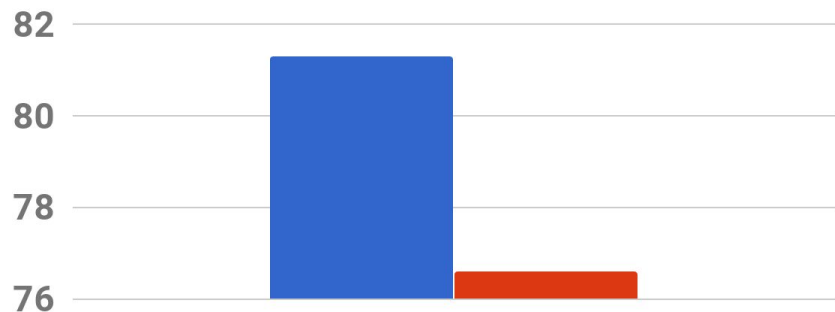
I ate pizza

I ate pizza

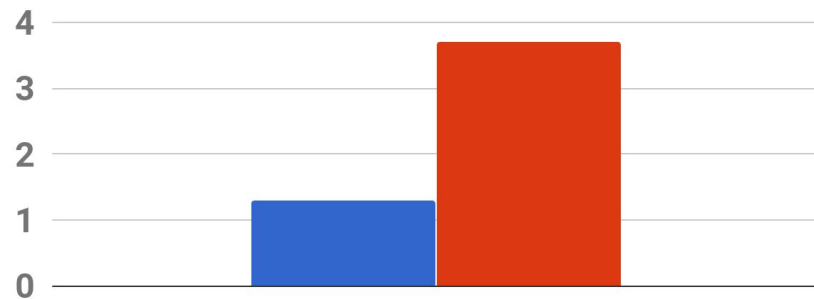# Or, "End-to-End Style" with Features from Tagger

# Pipeline: Higher Accuracy for Smaller Size

**Reordering Score**



**Model Size (MB)**



■ **Pipeline:** POS tags → Preordering

■ **End-to-end:** Preordering+Tagger features

Google

# Case Study 3: Language Identification

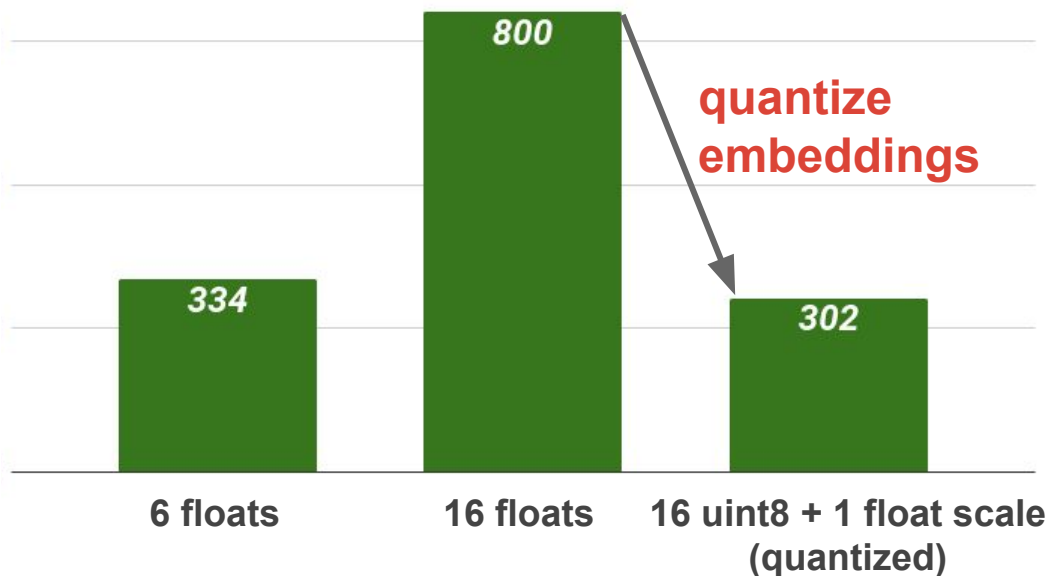# LangID: Post-hoc Quantization to Reduce Space

*Baldwin & Lui (2010)*
*.870 - .902*

F1:    .873        .880        .880

Model size (KB)



| 6 floats | 16 floats | 16 uint8 + 1 float scale (quantized) |

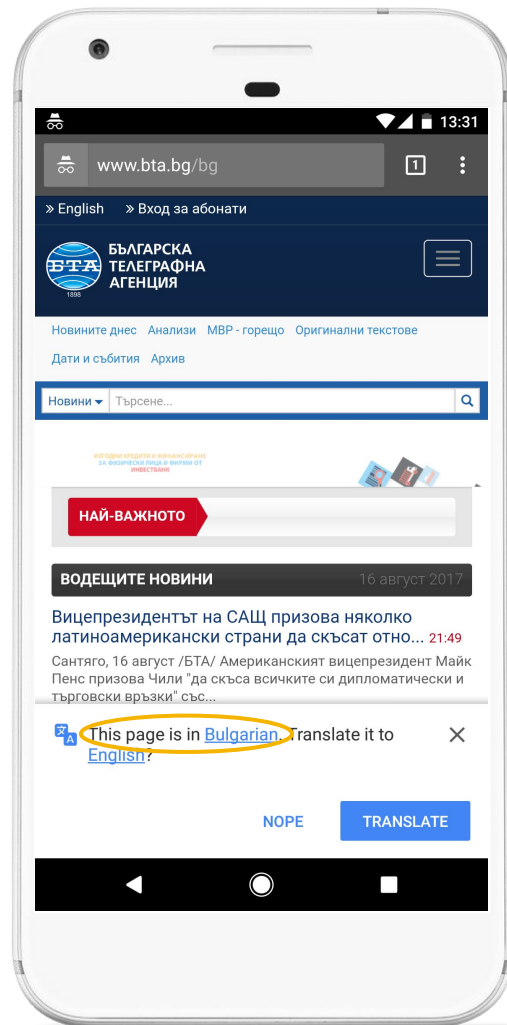**6 dimensions**        **16 dimensions**

quantize embeddings

800

334    302

Google

# That LangID model is essentially...

*Compact Language Detector v3 (CLD3)*

✓ runs inside all Google Chrome browsers

✓ code: github.com/google/CLD3

*Actual screenshot from 16 Aug 2017*

# Conclusion

Small (<= 2 MB) & fast (7k - 46k tokens/second) models with high accuracy in multiple tasks

Explicit intermediate representations & engineered features bring big accuracy gains in low-memory setting

Simple techniques → easy to include in standard practice

Google