

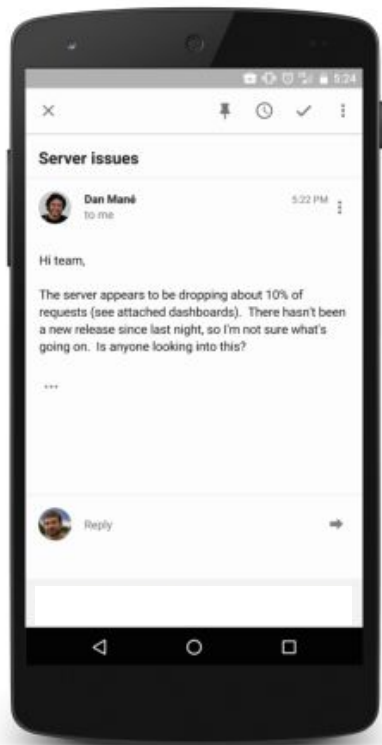


Deep Learning for Language Understanding (at Google Scale)

Anjuli Kannan

Software Engineer, Google Brain

Text is just a sequence of words



["hi", "team", "the",
"server", "appears", "to",
"be", "dropping", "about",
"10%", ...]

About me

- My team: Google Brain
 - "Make machines intelligent, improve people's lives."
 - Research + software + applications
 - g.co/brain
- My work is at boundary of research and applications
- Focus on natural language understanding

Neural network basics

Neural network

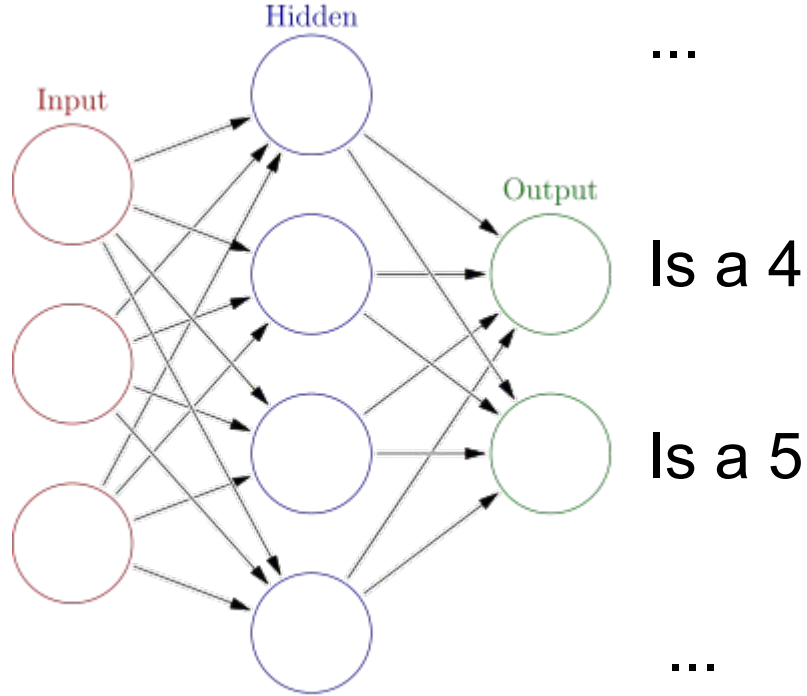
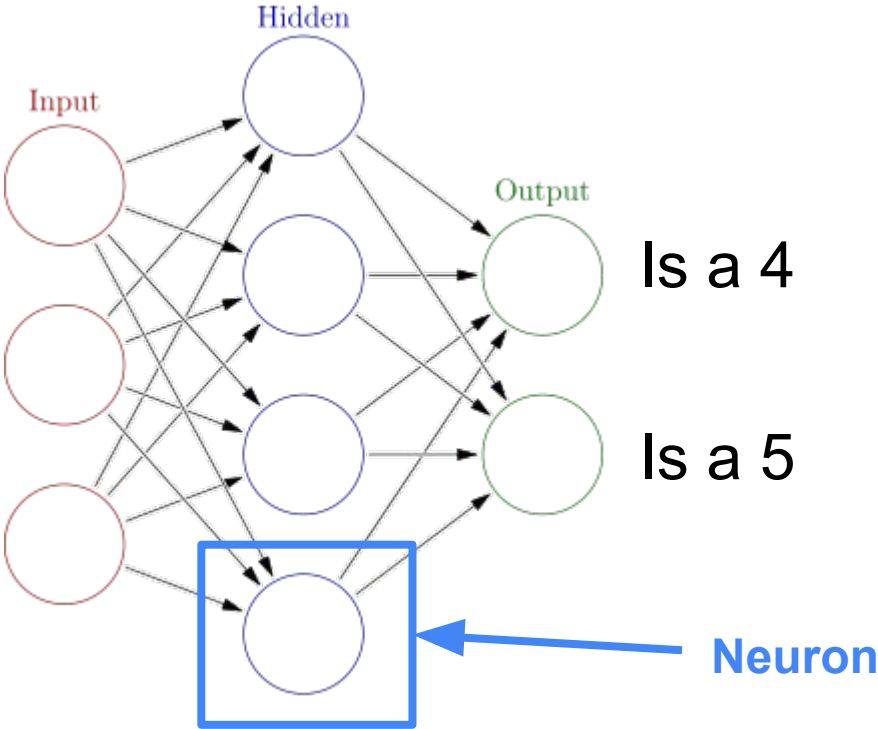


Image: Wikipedia

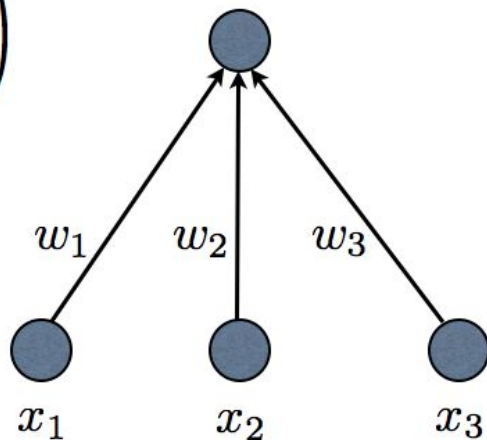
Neural network



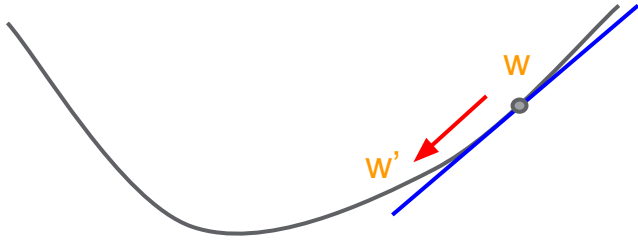
Basic building block is the neuron

- Different weights compute different functions

$$y_i = F \left(\sum_i w_i x_i \right)$$



Gradient descent



$$w' = w - \alpha \partial_w L(w)$$

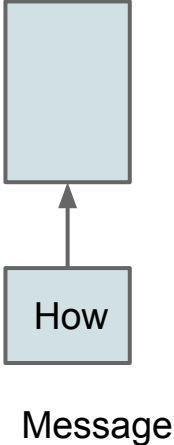
Learning Rate

Slide: Vincent Vanhoucke

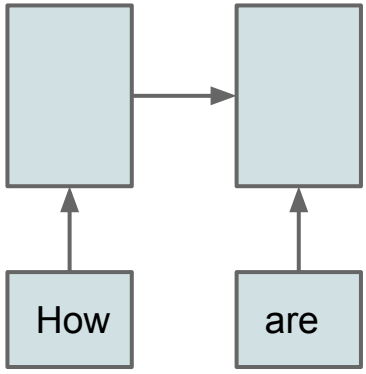
Recurrent neural networks

Recurrent neural networks can model sequences

Recurrent neural networks can model sequences

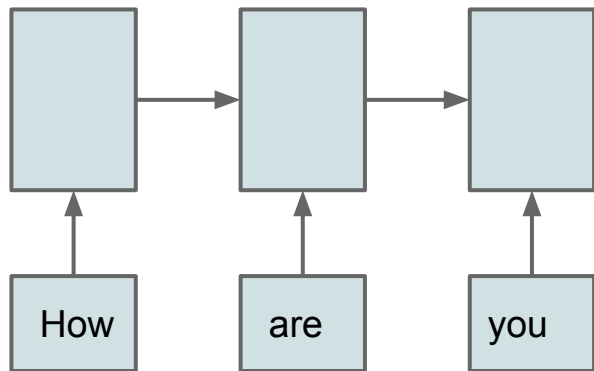


Recurrent neural networks can model sequences



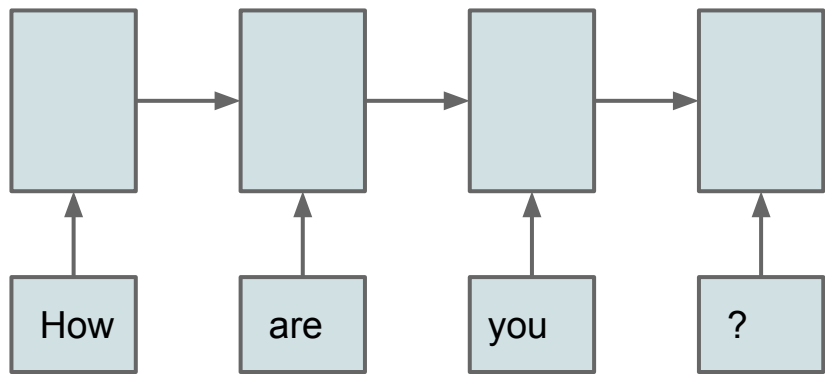
Message

Recurrent neural networks can model sequences



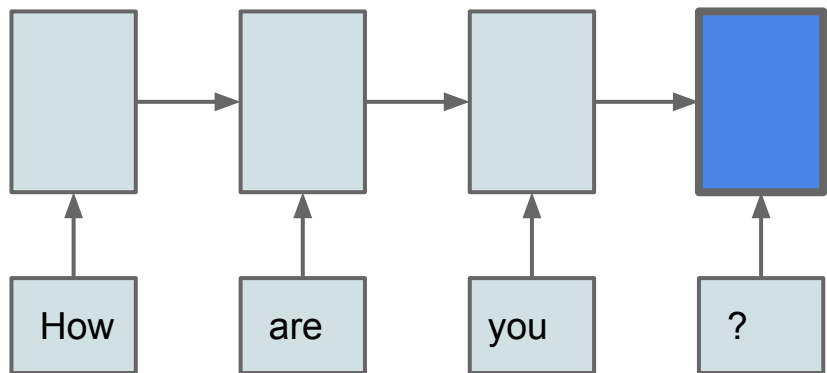
Message

Recurrent neural networks can model sequences



Message

Recurrent neural networks can model sequences



Message

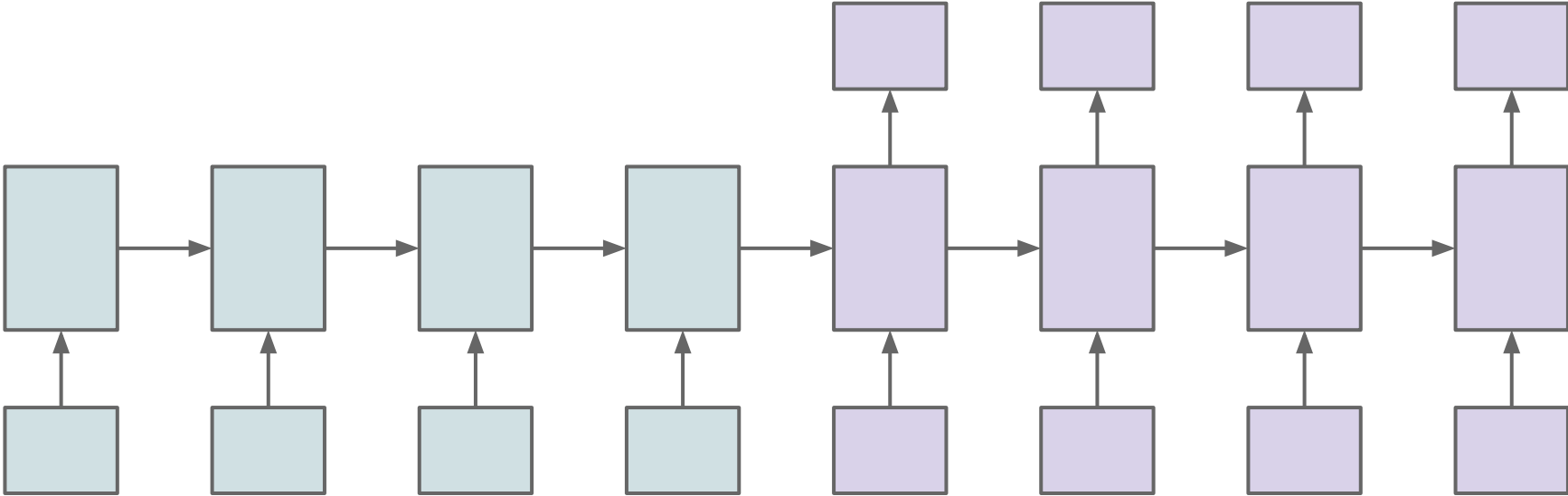
Internal state is a fixed length encoding of the message

Sequence-to-sequence models

Suppose we want to generate email replies

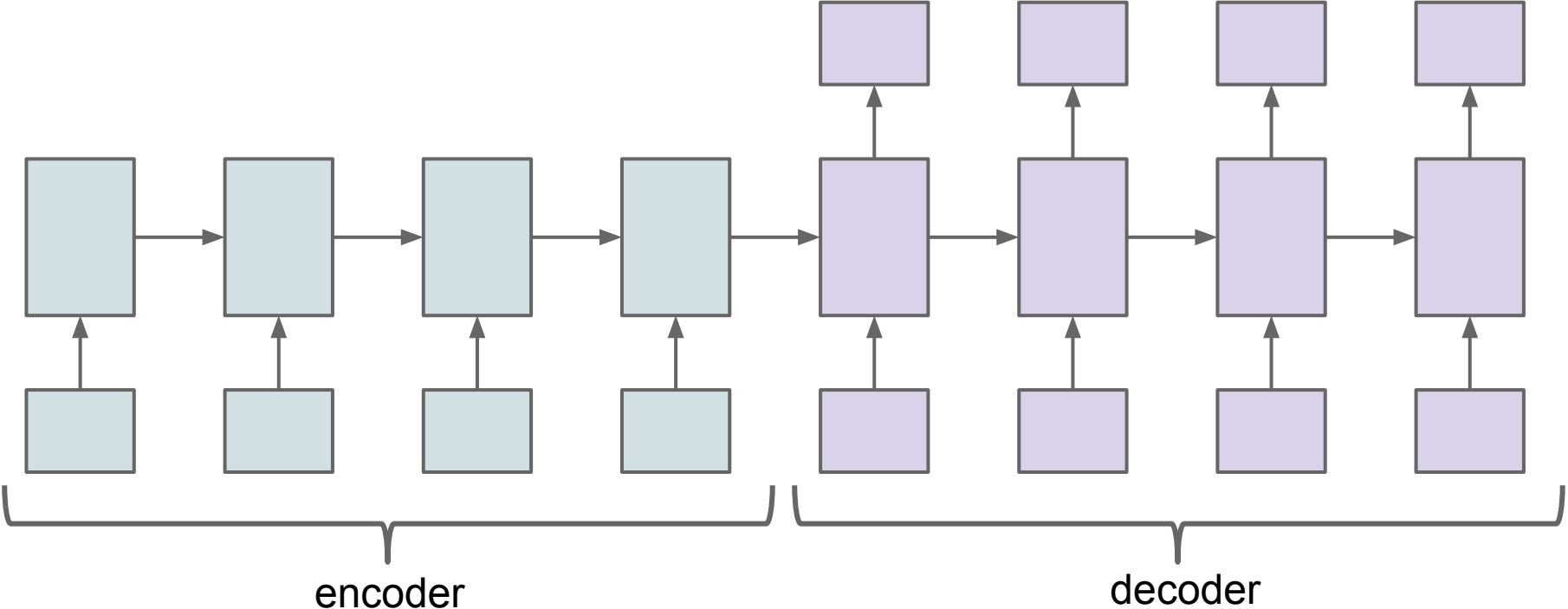


Sequence-to-sequence model

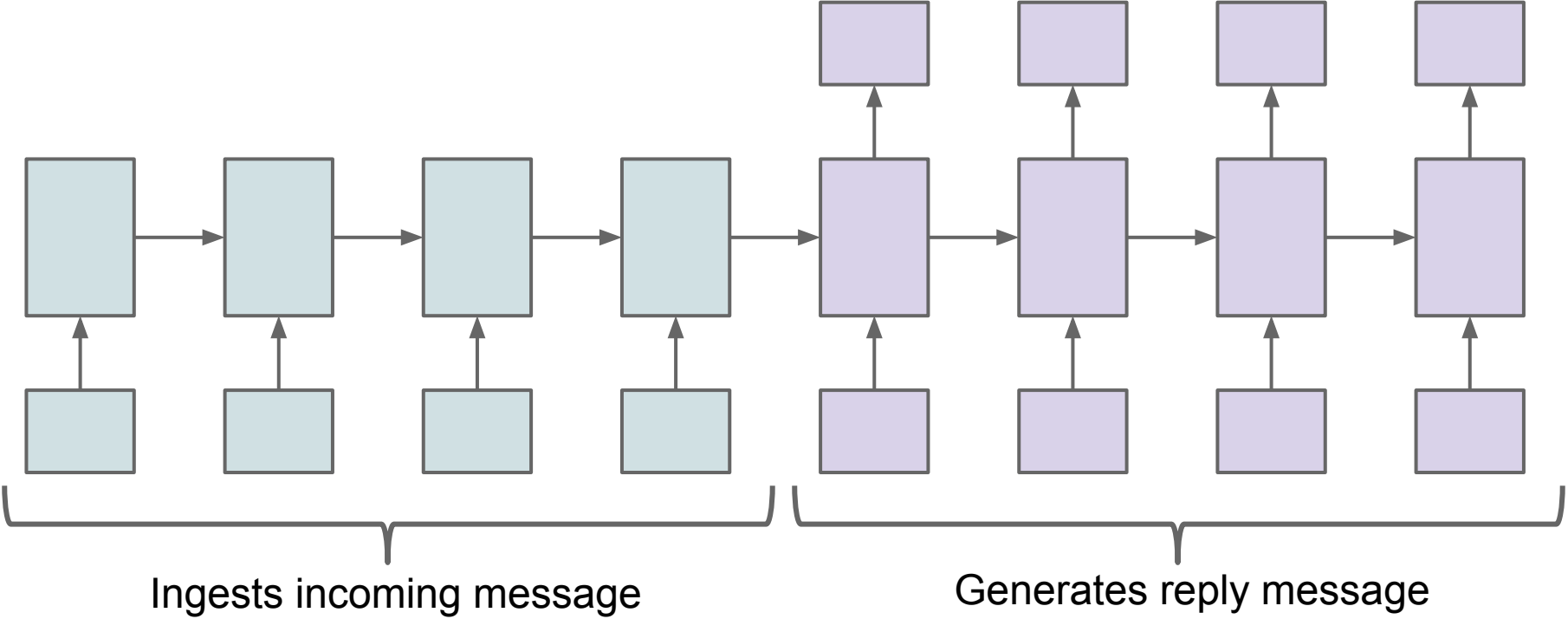


Sutskever et al, NIPS 2014

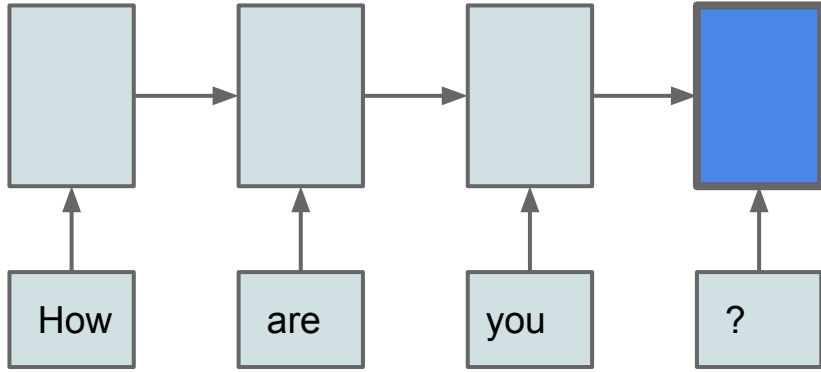
Sequence-to-sequence model



Sequence-to-sequence model



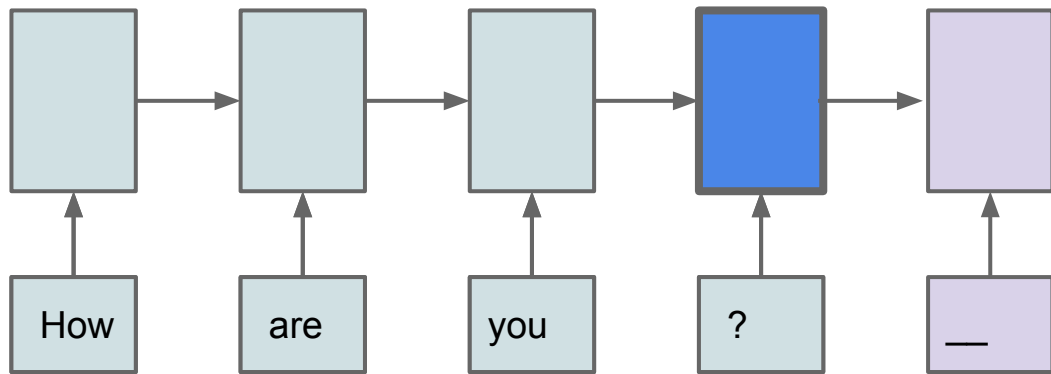
Encoder ingests the incoming message



Message

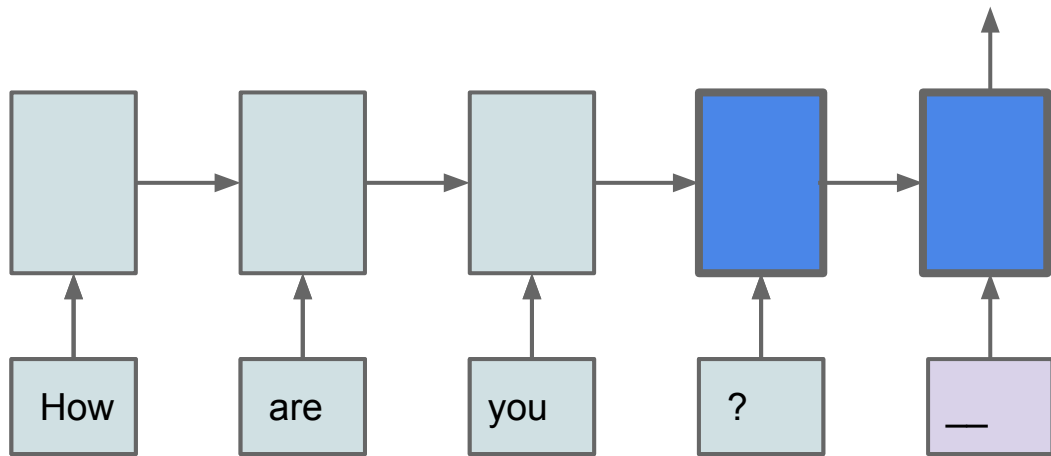
Internal state is a fixed length encoding of the message

Decoder is initialized with final state of encoder



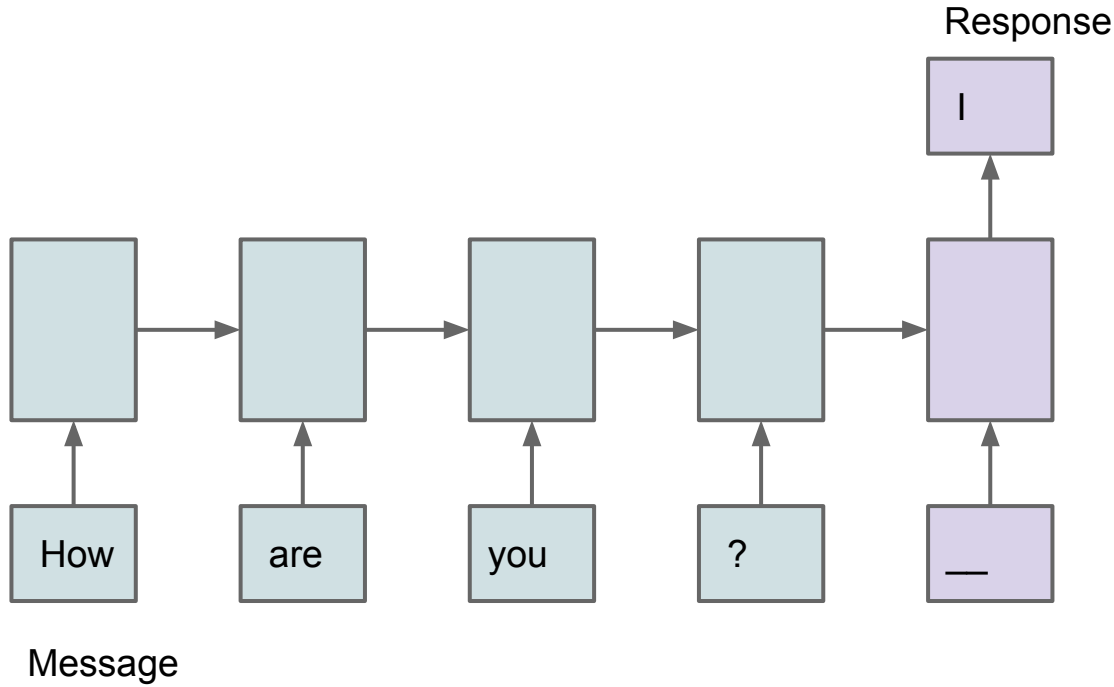
Message

Decoder is initialized with final state of encoder

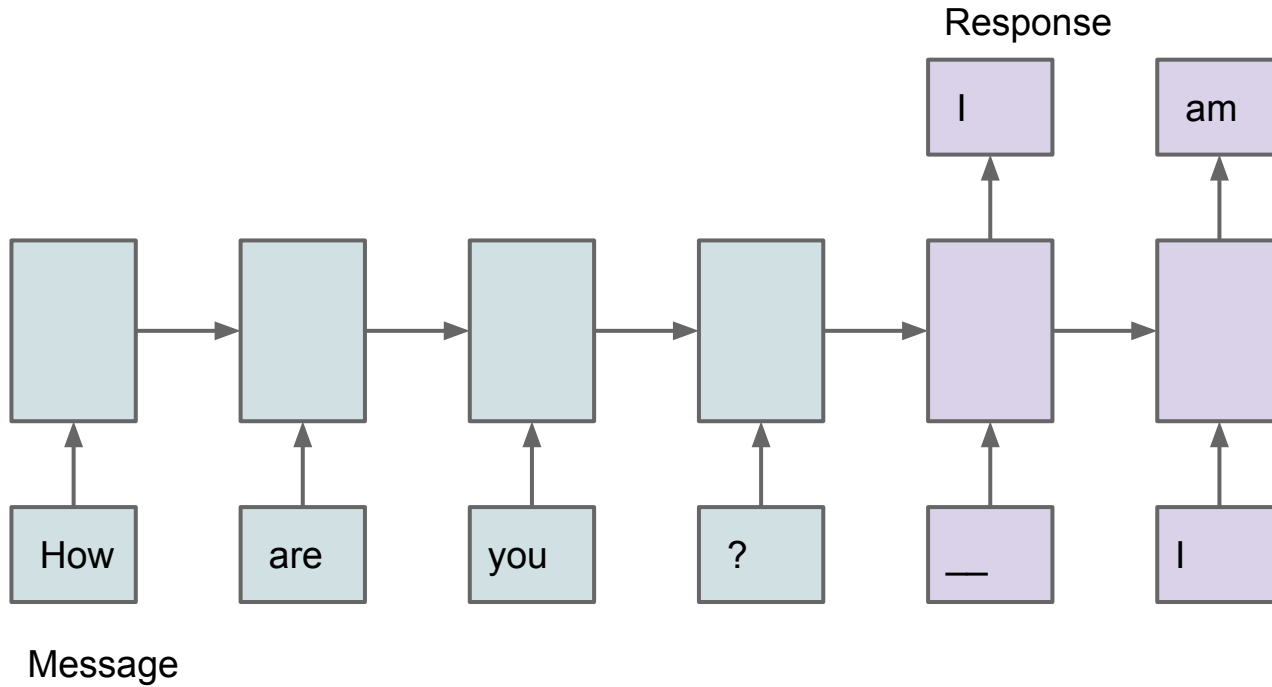


Message

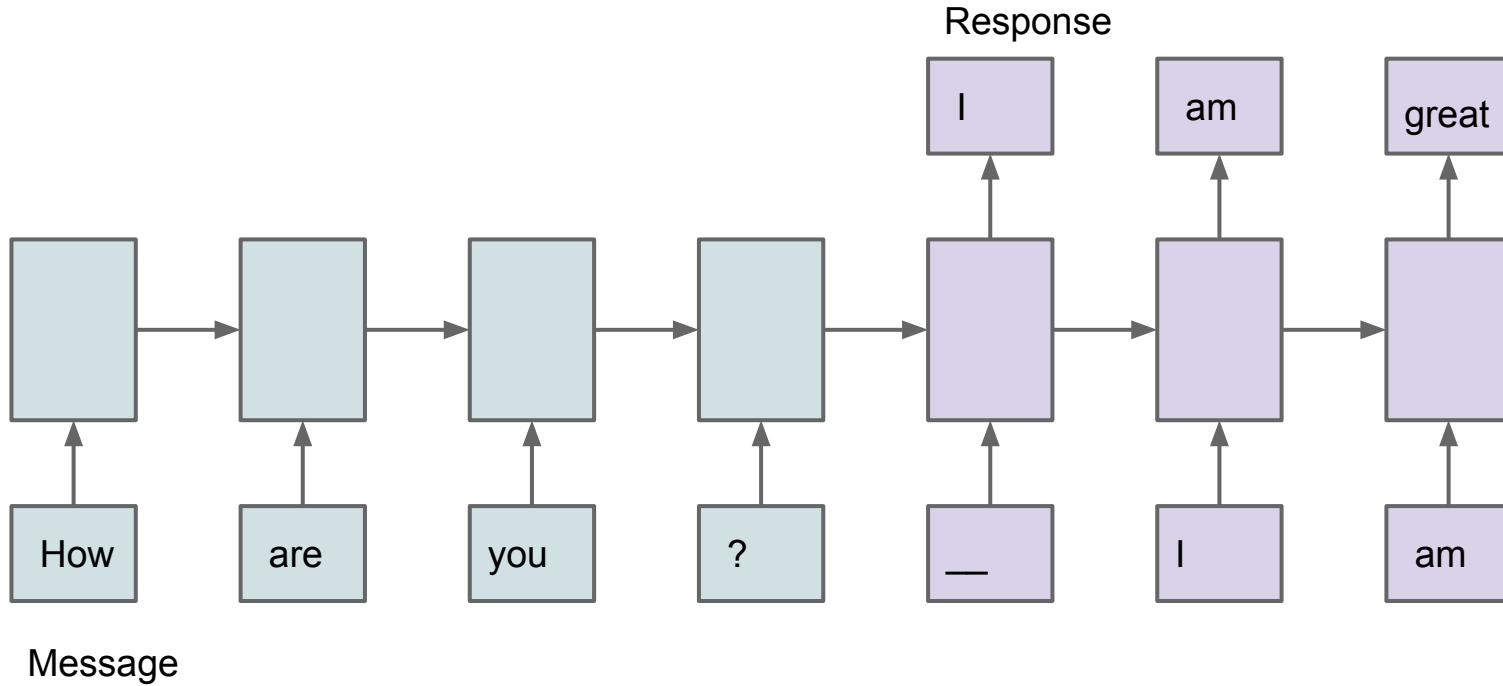
Decoder predicts next word



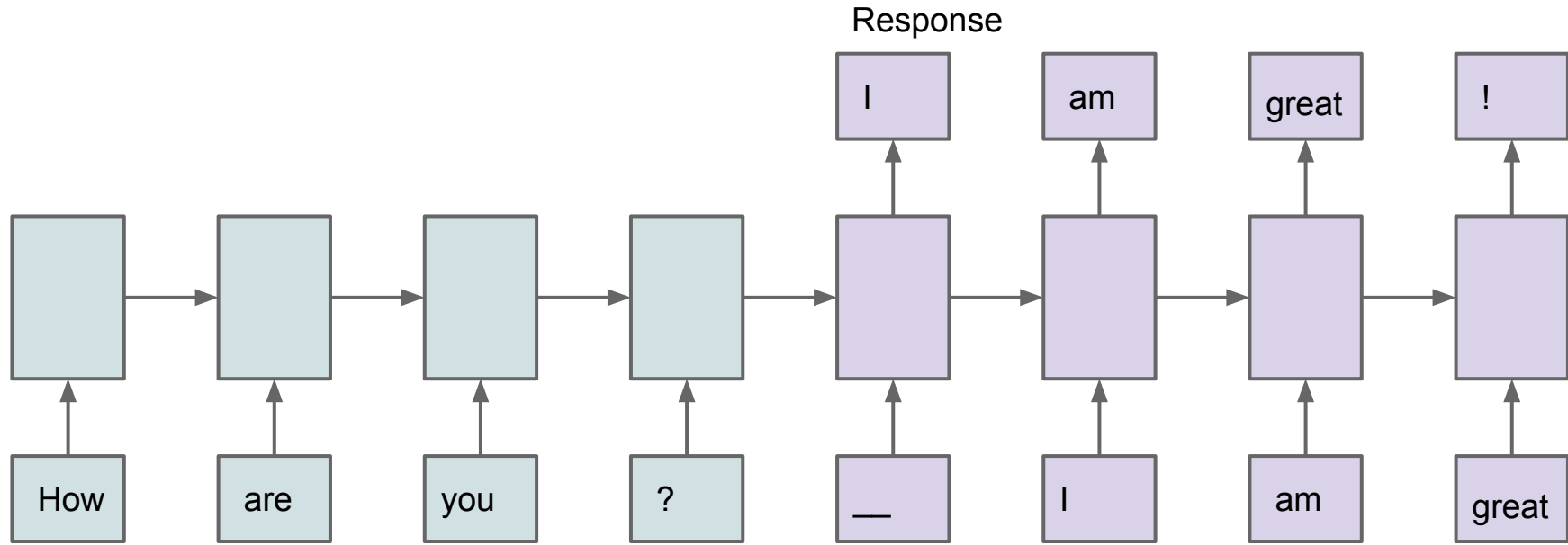
Decoder predicts next word



Decoder predicts next word



Decoder predicts next word



Message

Response

Vinyals & Le, ICML DL 2015
Kannan et al, KDD 2016

What the model can do

message = Can you do Tuesday or Wednesday?

responses :

I can do Tuesday

I can do Wednesday

I can do Wednesday .

I can do Tuesday .

Tuesday works for me .

Wednesday works for me .

Tuesday works for me

Tuesday works .

Tuesday works for me !

Yes I can

What the model can do

message = I feel so gross this morning. I think I ate something bad last night. Feel like I'm going to barf.

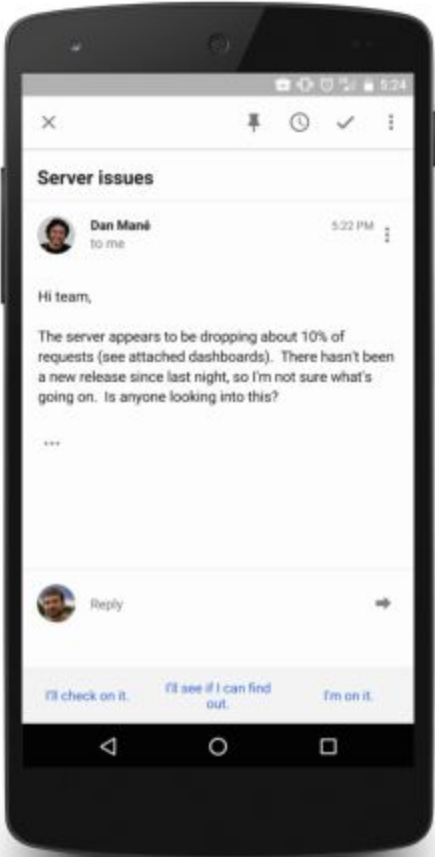
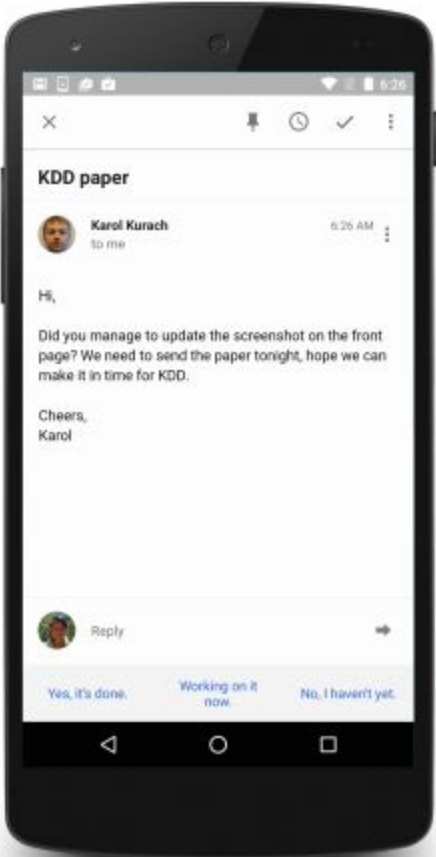
responses:

Feel better !
What 's wrong ?
I 'm sorry .
I 'm sorry :(
I 'm sorry. :(
Feel better .
Oh no !
What did you eat ?
I 'm sorry
That sucks .

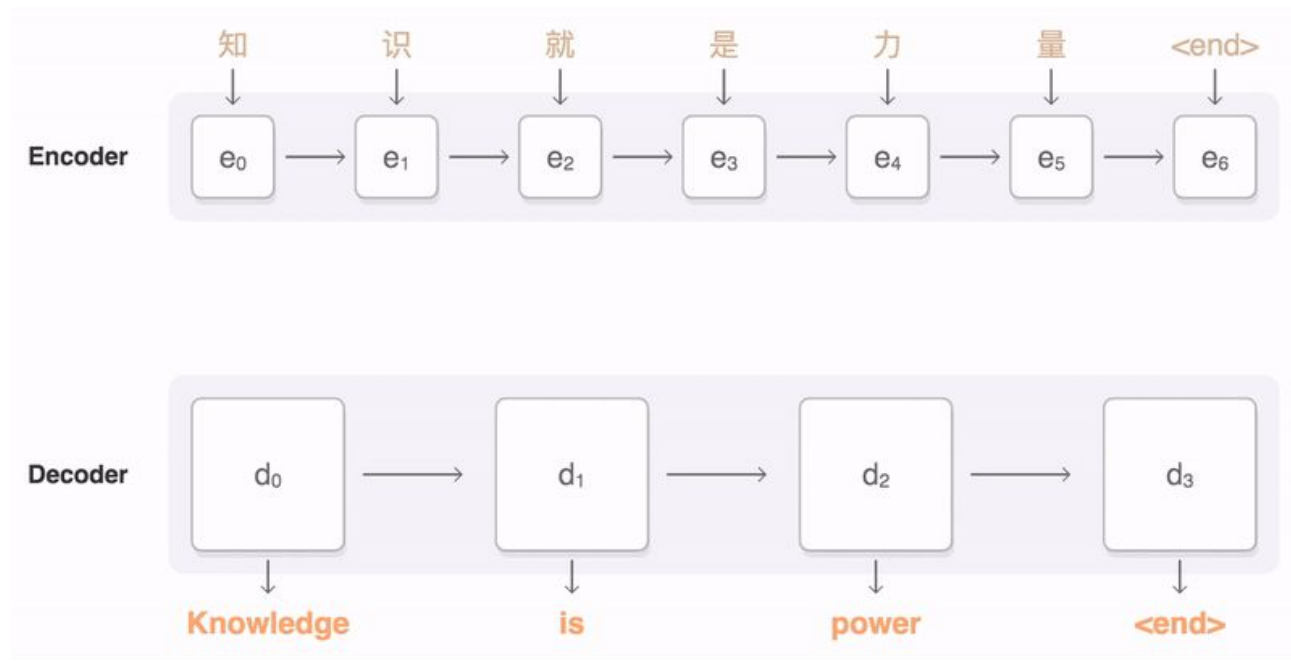
Summary

- Neural networks learn **feature representations** from raw data
- **Recurrent neural networks** have statefulness which allows them to model sequences of data such as text
- The **sequence-to-sequence** model contains two recurrent neural networks: one to encode an input sequence and one to generate an output sequence

Smartreply



Google Translate



Authors

Yonghui Wu, Mike Schuster, Zhirfeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, Jeffrey Dean

Research: Speech recognition

STATE-OF-THE-ART SPEECH RECOGNITION WITH SEQUENCE-TO-SEQUENCE MODELS

*Chung-Cheng Chiu, Tara N. Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen,
Zhifeng Chen, Anjuli Kannan, Ron J. Weiss, Kanishka Rao, Ekaterina Gonina,
Navdeep Jaitly, Bo Li, Jan Chorowski, Michiel Bacchiani*

Research: Electronic health records



1

Health systems collect and store electronic health records in various formats in databases.

Scalable and accurate deep learning for electronic health records

Alvin Rajkomar^{*1,2}, Eyal Oren^{*1}, Kai Chen¹, Andrew M. Dai¹, Nissan Hajaj¹, Peter J. Liu¹, Xiaobing Liu¹, Mimi Sun¹, Patrik Sundberg¹, Hector Yee¹, Kun Zhang¹, Yi Zhang¹, Gavin E. Duggan¹, Gerardo Flores¹, Michaela Hardt¹, Jamie Irvine¹, Quoc Le¹, Kurt Litsch¹, Jake Marcus¹, Alexander Mossin¹, Justin Tansuwan¹, De Wang¹, James Wexler¹, Jimbo Wilson¹, Dana Ludwig², Samuel L. Volchenboum⁴, Katherine Chou¹, Michael Pearson¹, Srinivasan Madabushi¹, Nigam H. Shah³, Atul J. Butte², Michael Howell¹, Claire Cui¹, Greg Corrado¹, and Jeff Dean¹

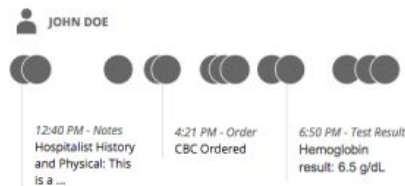
¹Google Inc, Mountain View, California

²University of California, San Francisco, San Francisco, California

³Stanford University, Stanford, California

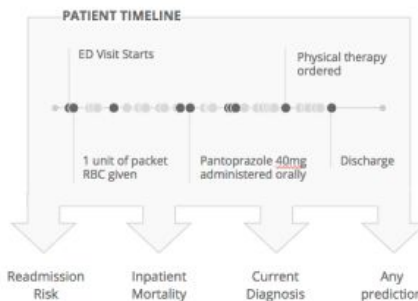
⁴University of Chicago Medicine, Chicago, Illinois

January 2018



2

All available data for each patient is converted to events recorded in containers based on the Fast Healthcare Interoperability Resource (FHIR) specification.



3

The FHIR resources are placed in temporal order, depicting all events recorded in the EHR (i.e. timeline). The deep learning model uses this full history to make each prediction.

What's next?



Resources

- All tensorflow tutorials:
<https://www.tensorflow.org/versions/master/tutorials/index.html>
- Sequence-to-sequence tutorial (machine translation):
<https://www.tensorflow.org/versions/master/tutorials/seq2seq>
- Chris Olah's blog: <http://colah.github.io/>

Thank you!