

---

# Neural Machine Translation Decoding

Philipp Koehn

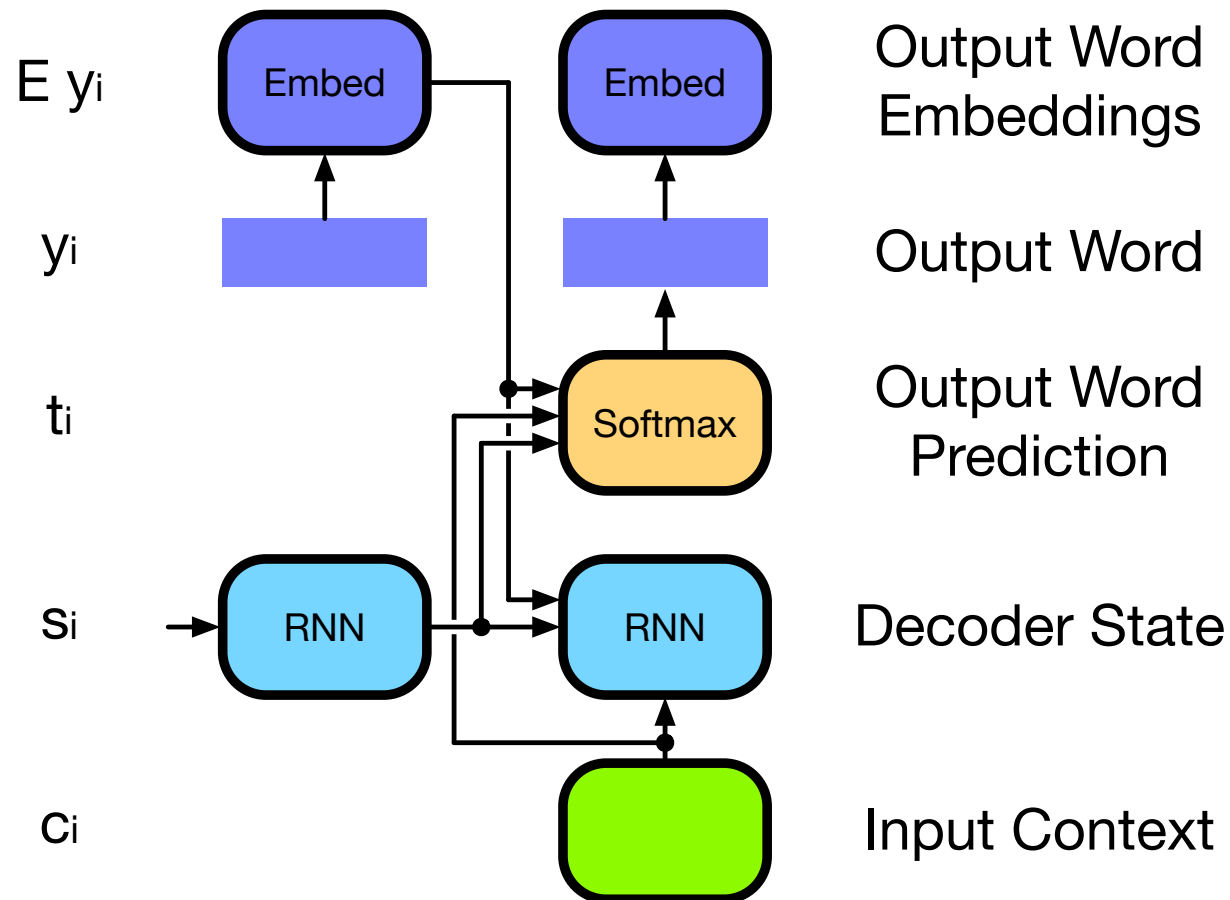


- Given a trained model  
... we now want to translate test sentences
- We only need execute the "forward" step in the computation graph

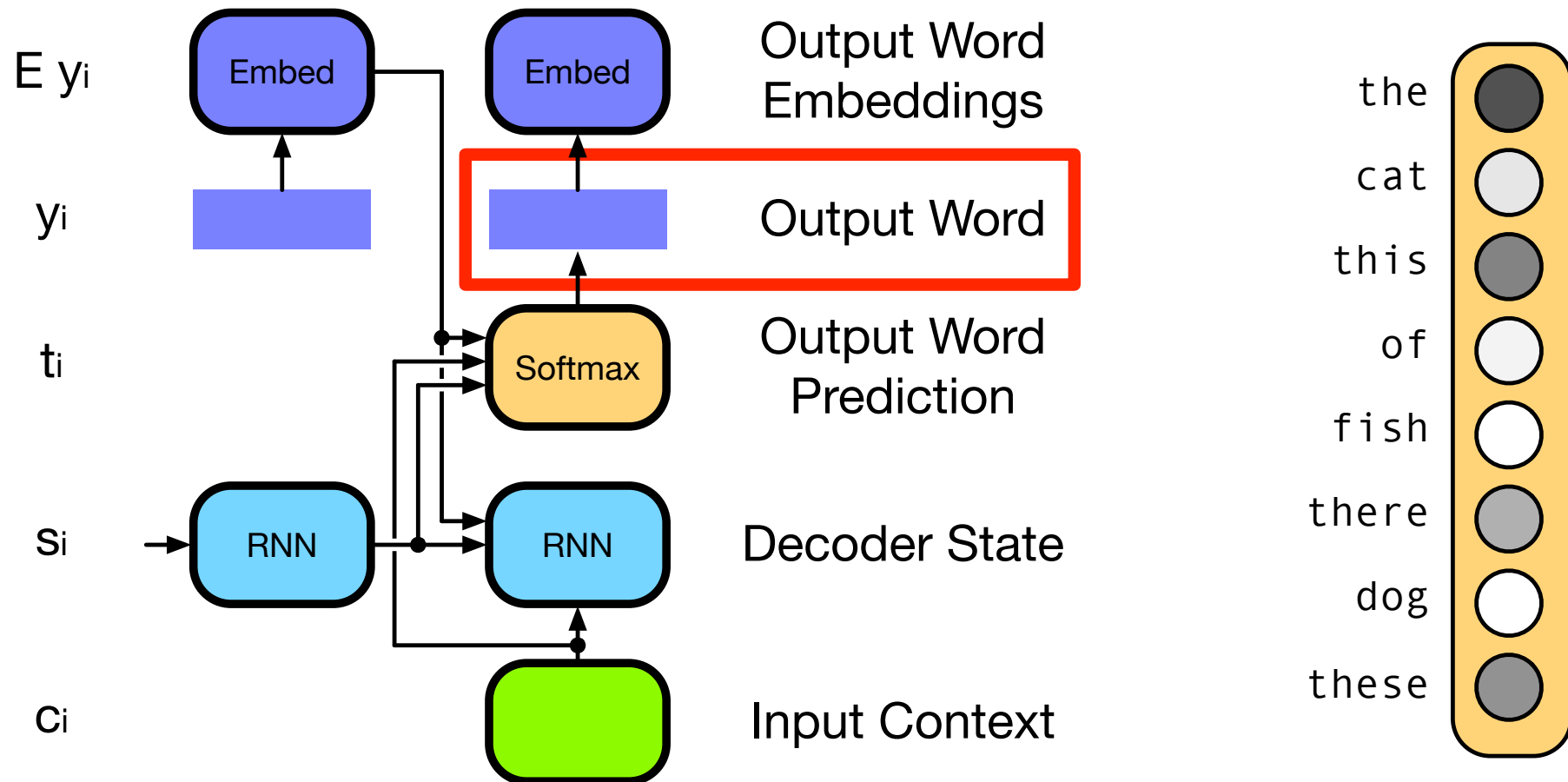
# Word Prediction



2



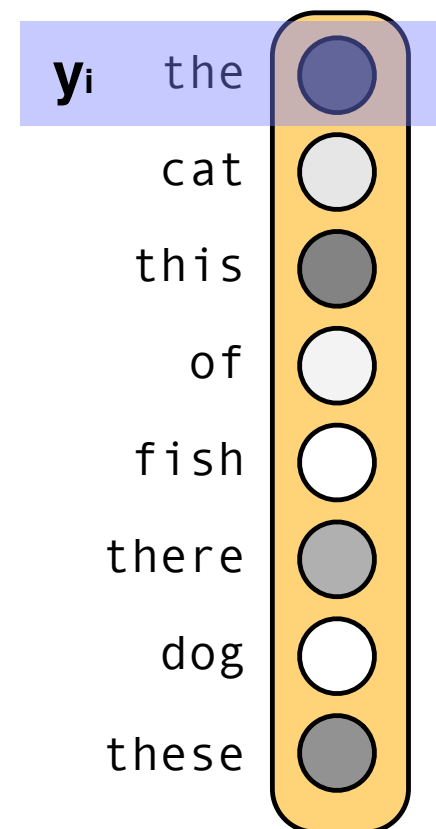
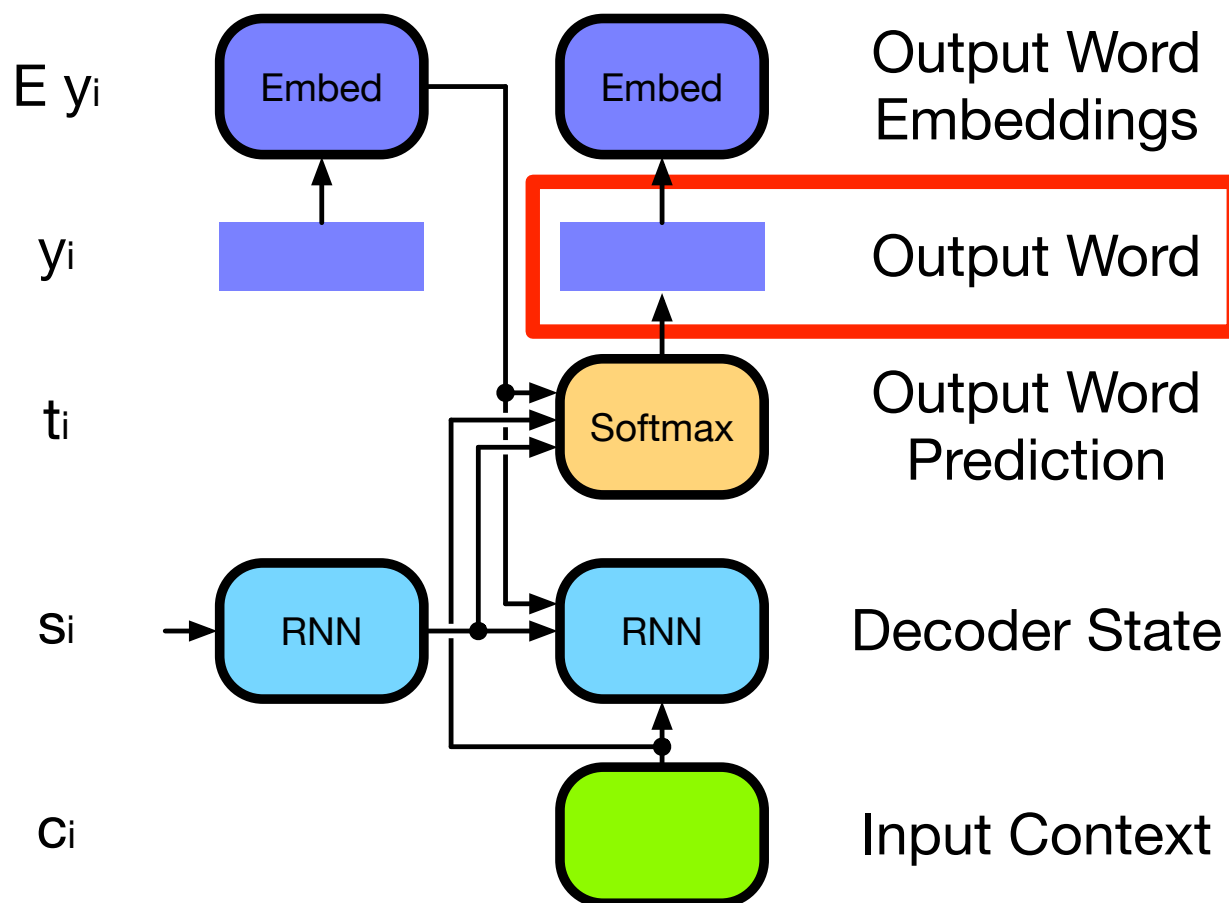
# Selected Word



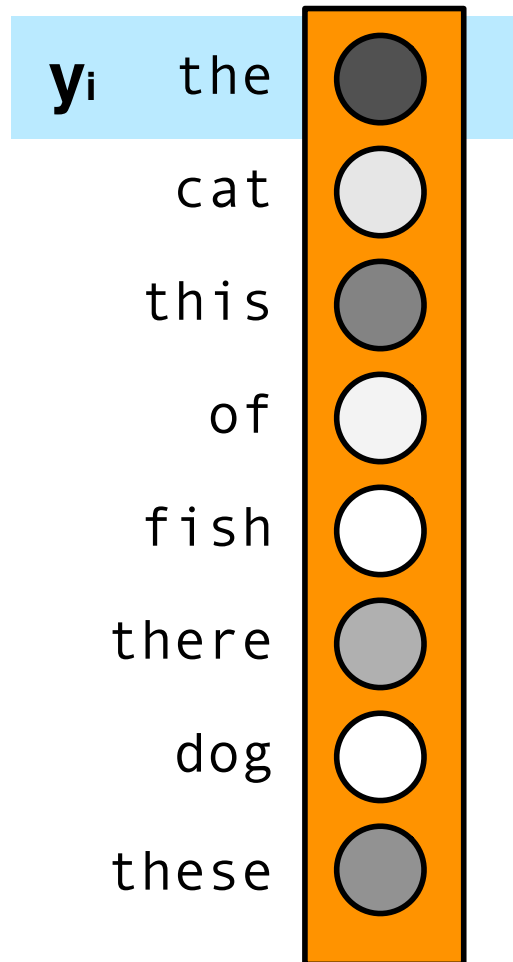
# Embedding



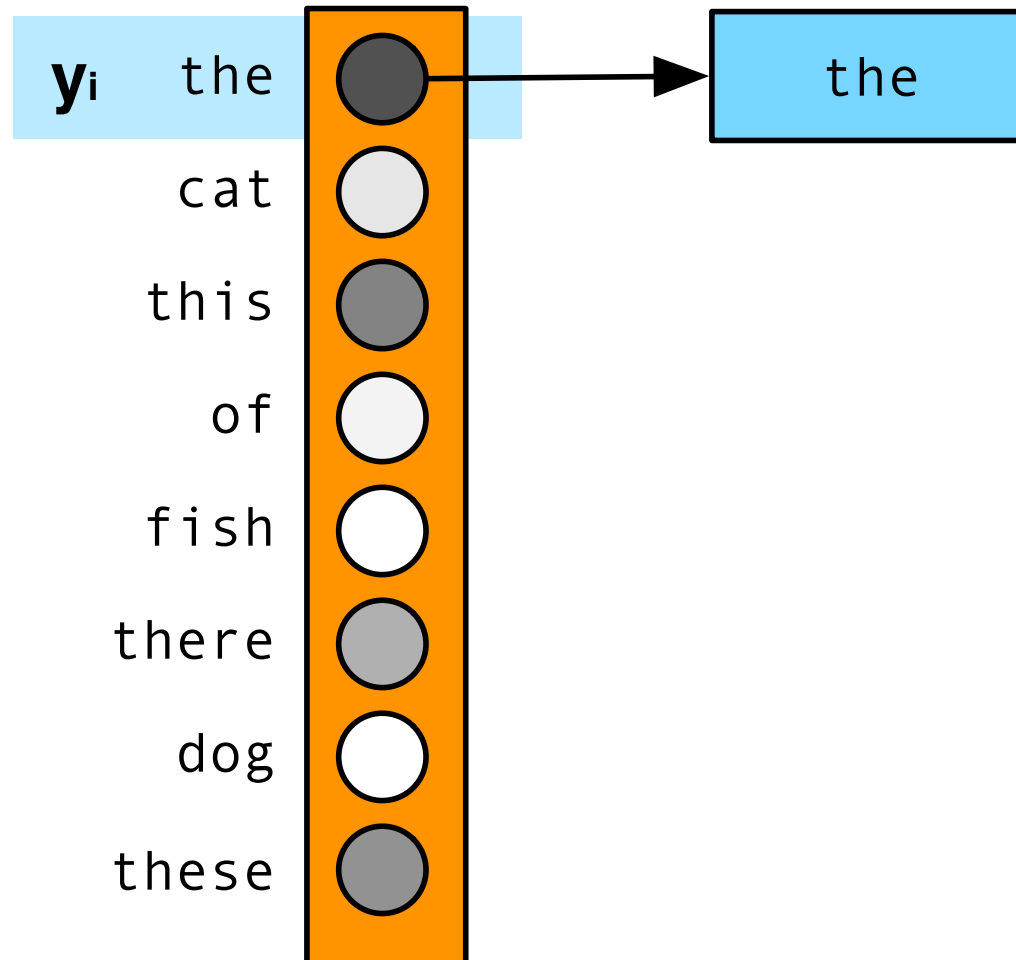
4



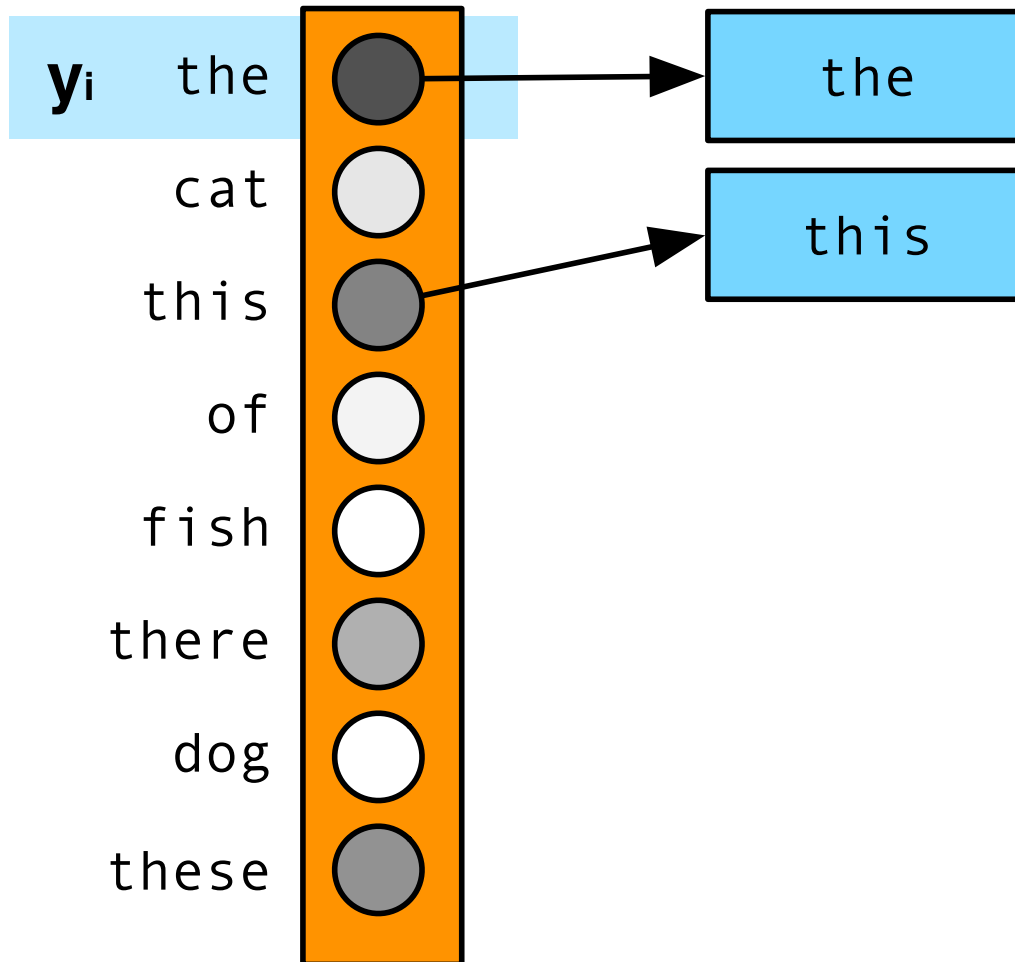
# Distribution of Word Predictions



# Select Best Word

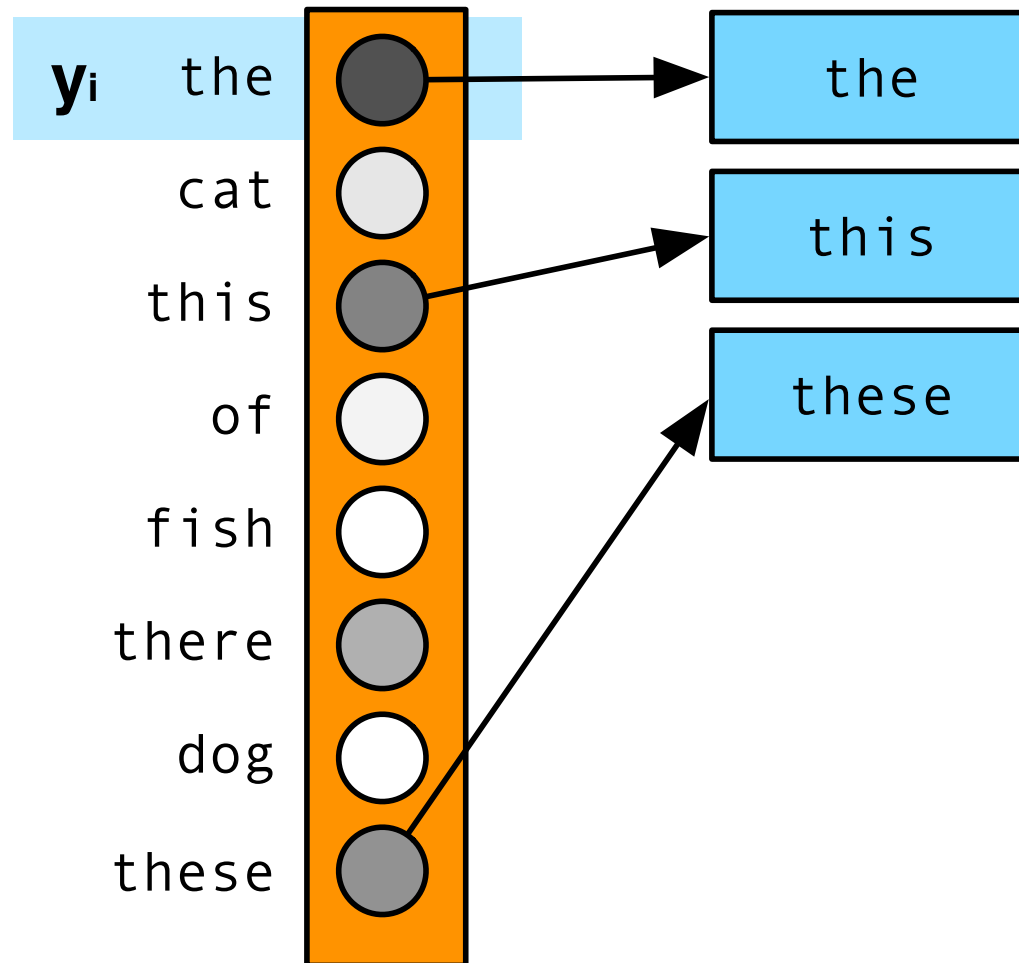


# Select Second Best Word





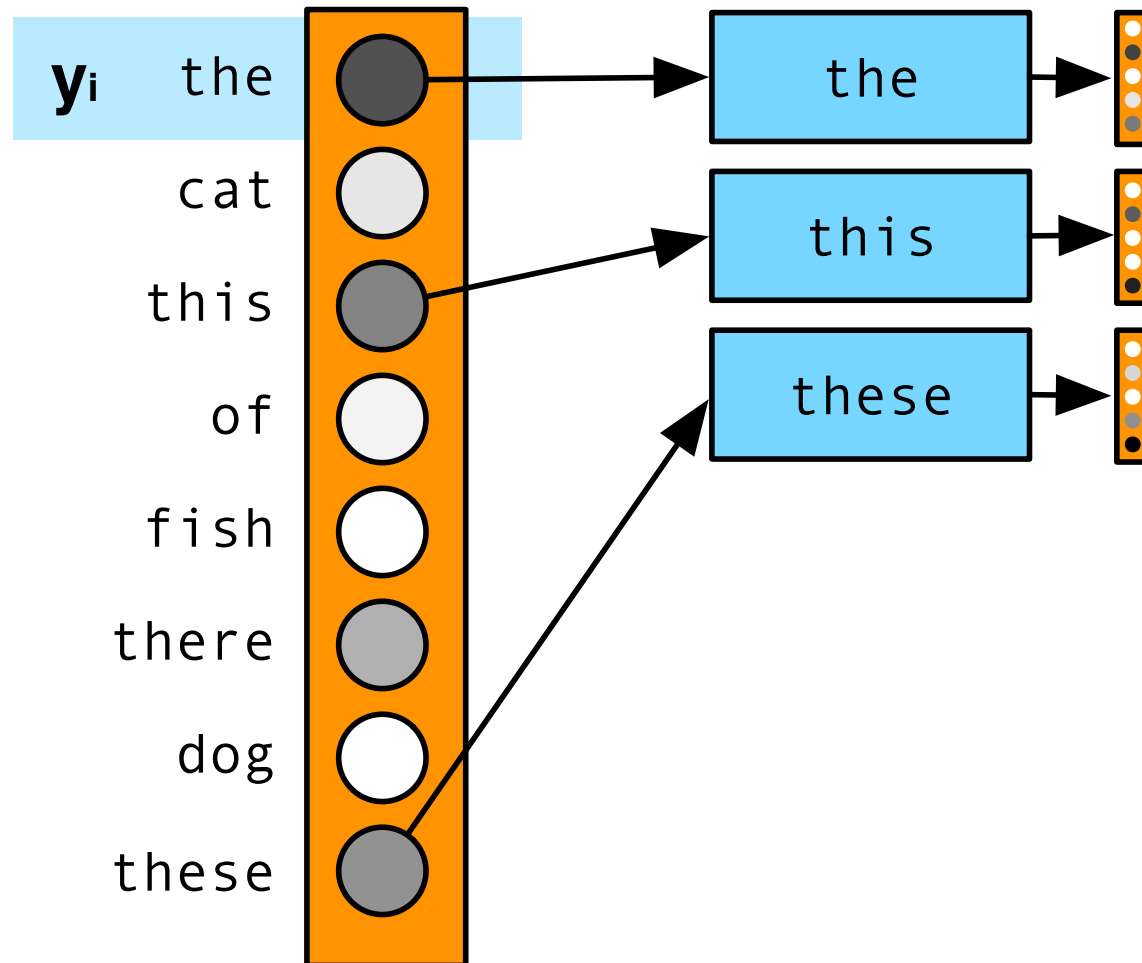
# Select Third Best Word



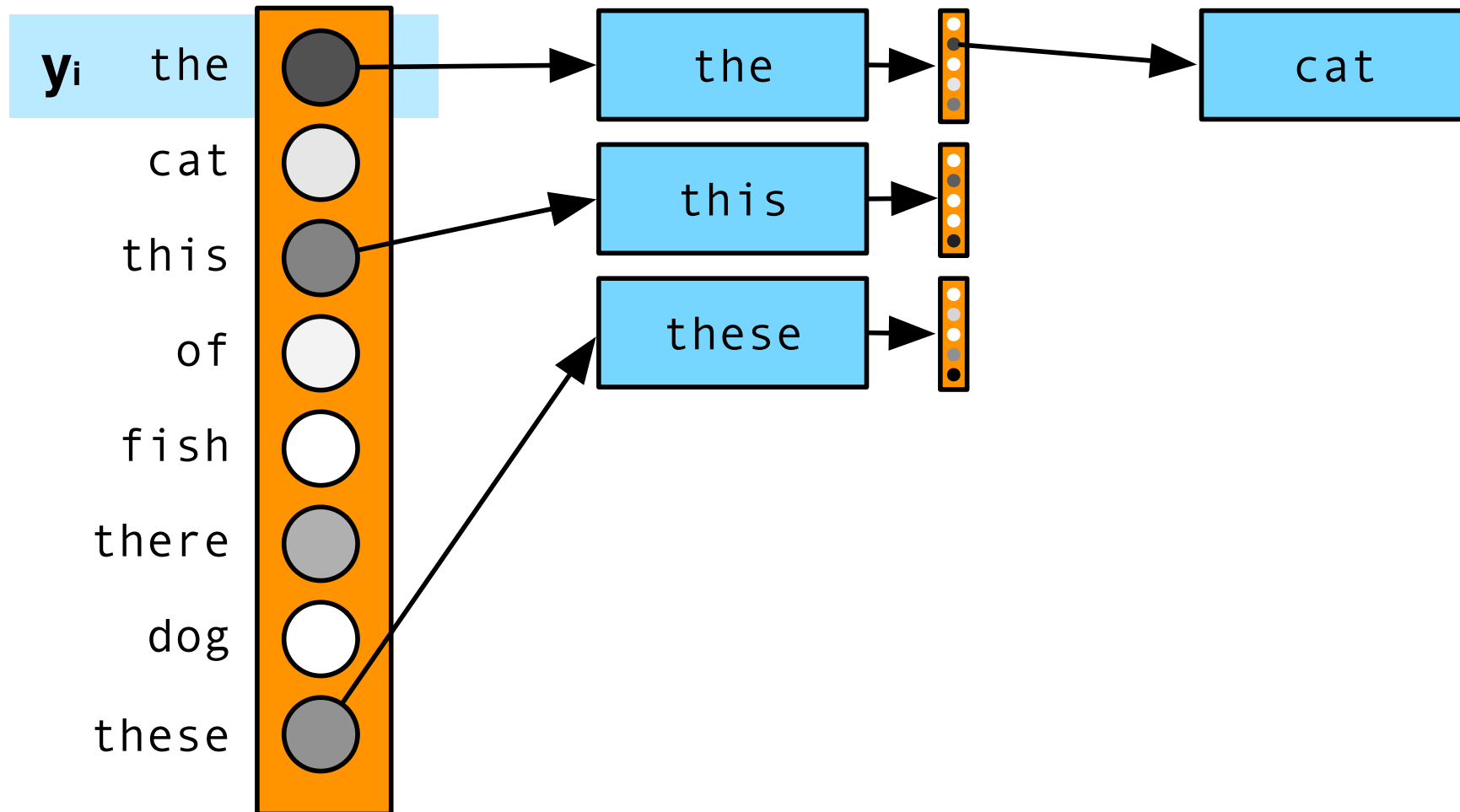
# Use Selected Word for Next Predictions



9

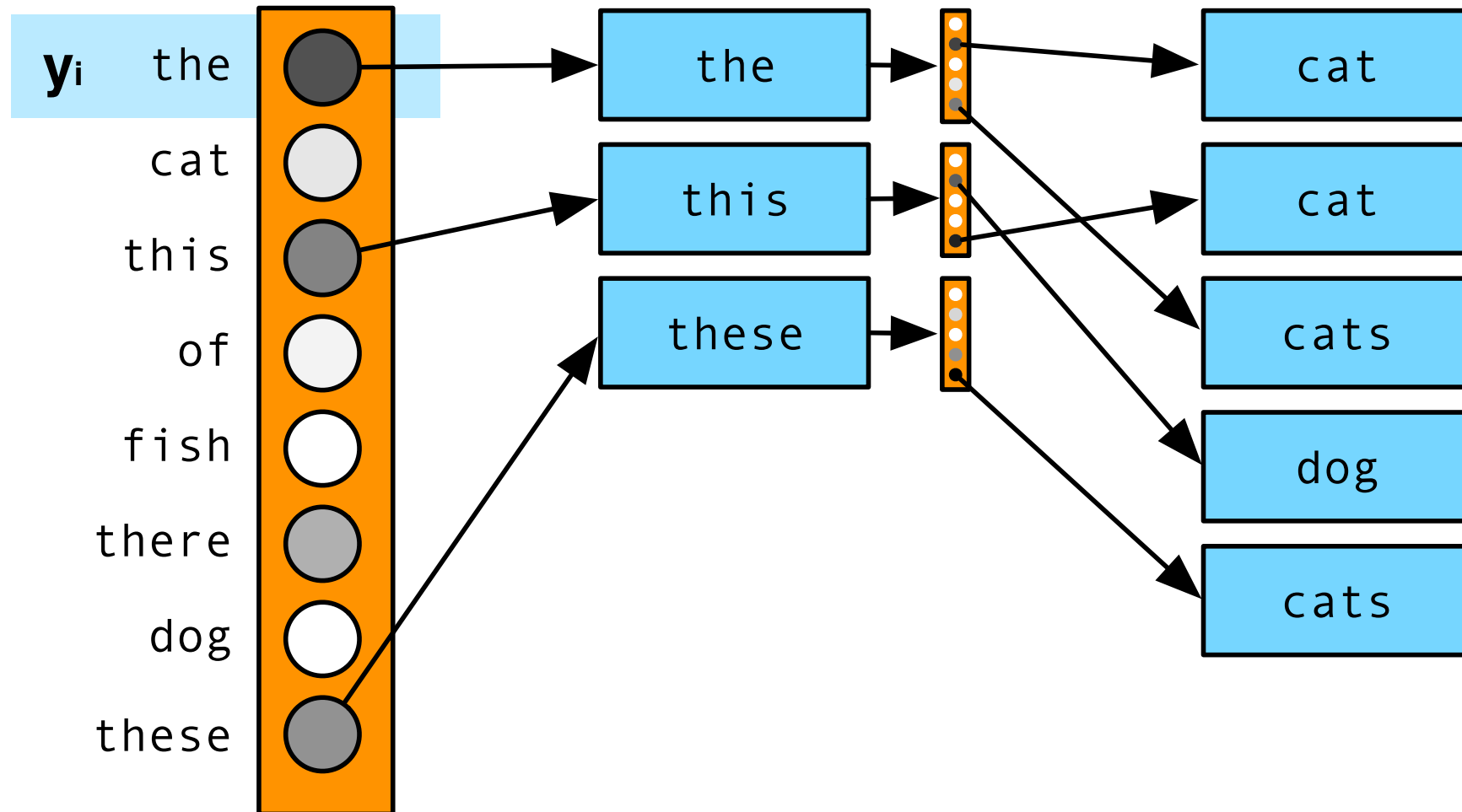


# Select Best Continuation



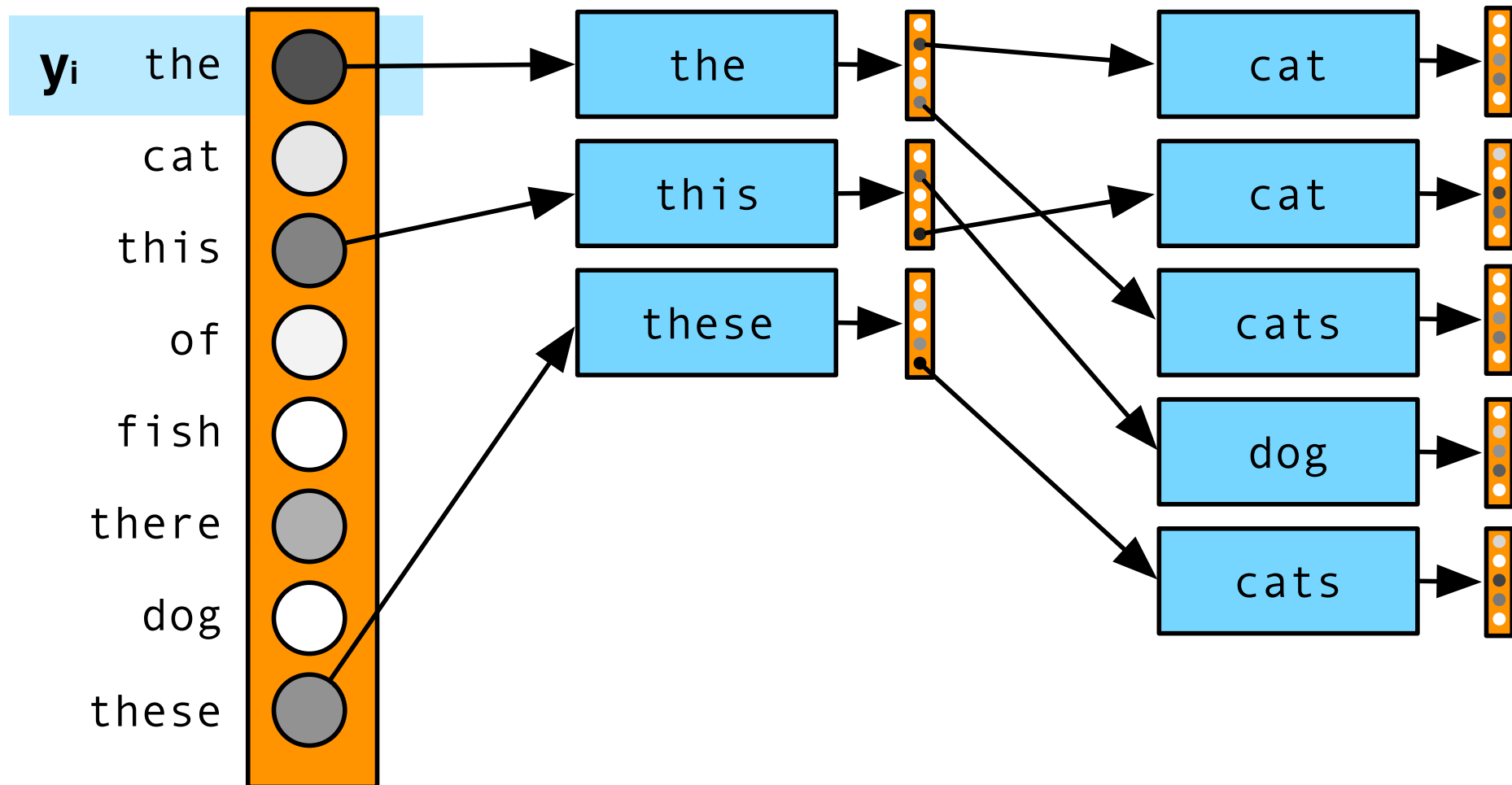
# Select Next Best Continuations

11

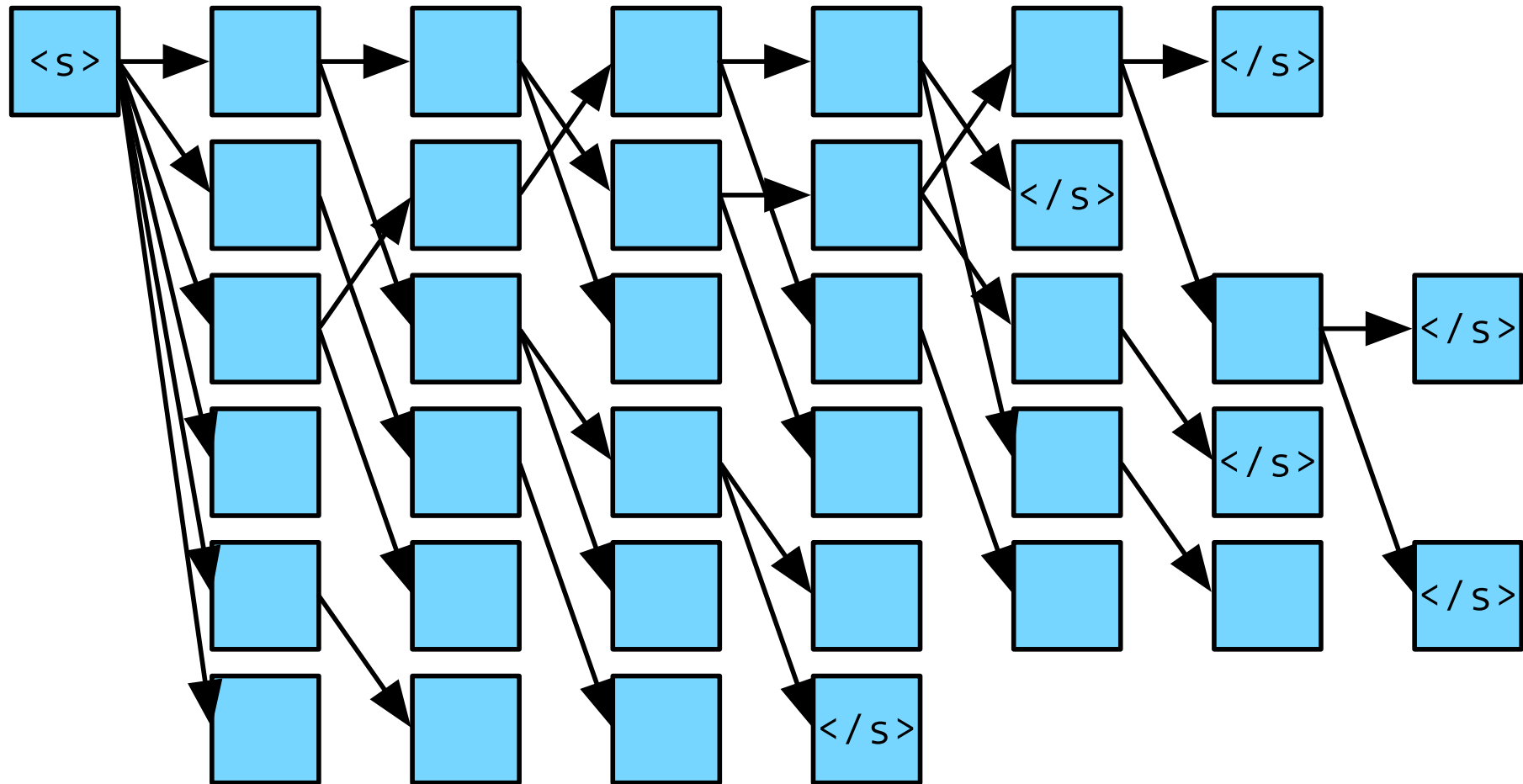


# Continue...

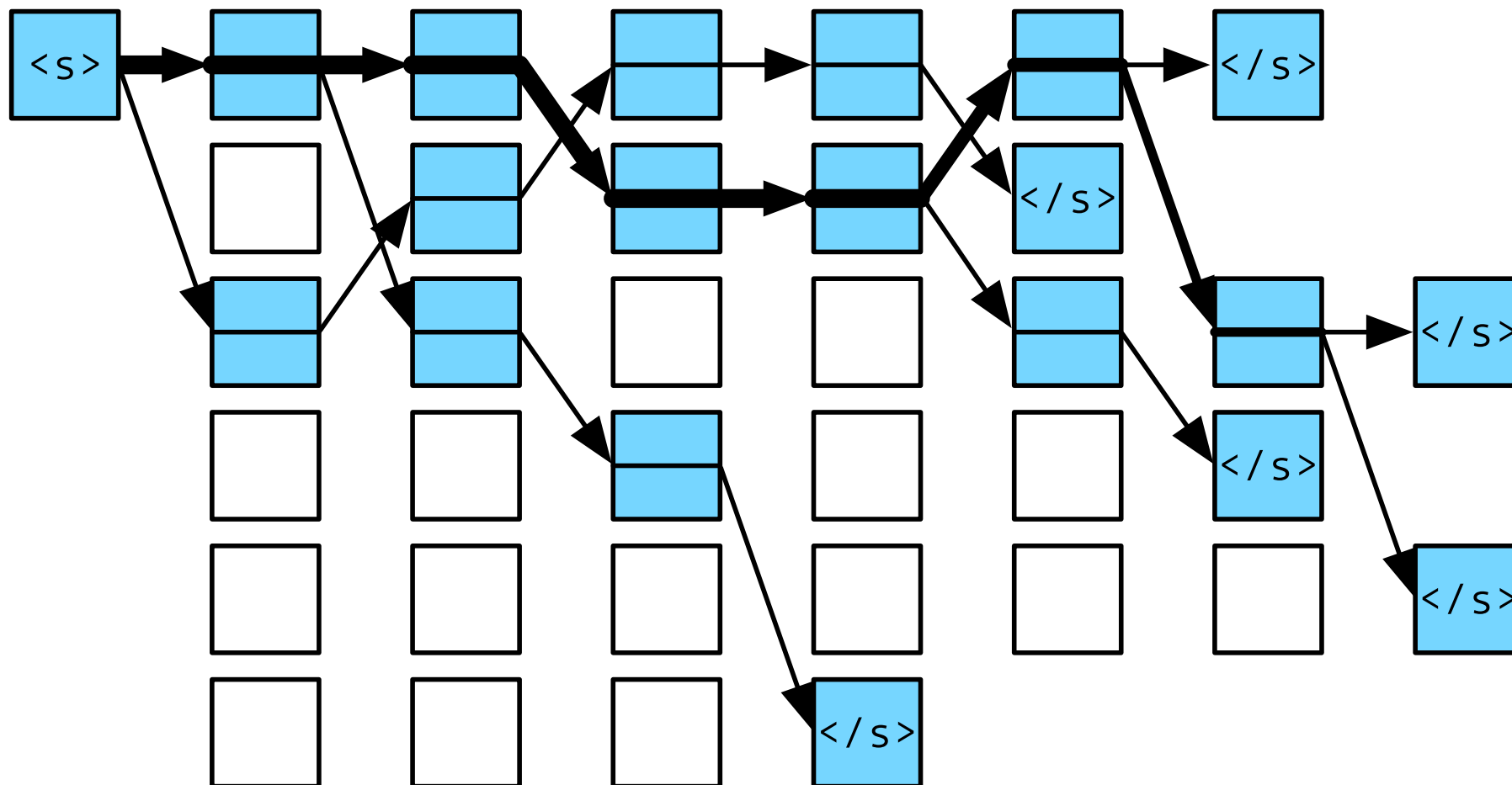
12



# Beam Search



## Best Paths



# Beam Search Details

- Normalize score by length
- No recombination (paths cannot be merged)



# Output Word Predictions

**Input Sentence:** *ich glaube aber auch , er ist clever genug um seine Aussagen vage genug zu halten , so dass sie auf verschiedene Art und Weise interpretiert werden können .*

Best		Alternatives
<b>but</b>	(42.1%)	<i>however (25.3%), I (20.4%), yet (1.9%), and (0.8%), nor (0.8%), ...</i>
<b>I</b>	(80.4%)	<i>also (6.0%), , (4.7%), it (1.2%), in (0.7%), nor (0.5%), he (0.4%), ...</i>
<b>also</b>	(85.2%)	<i>think (4.2%), do (3.1%), believe (2.9%), , (0.8%), too (0.5%), ...</i>
<b>believe</b>	(68.4%)	<i>think (28.6%), feel (1.6%), do (0.8%), ...</i>
<b>he</b>	(90.4%)	<i>that (6.7%), it (2.2%), him (0.2%), ...</i>
<b>is</b>	(74.7%)	<i>'s (24.4%), has (0.3%), was (0.1%), ...</i>
<b>clever</b>	(99.1%)	<i>smart (0.6%), ...</i>
<b>enough</b>	(99.9%)	
<b>to</b>	(95.5%)	<i>about (1.2%), for (1.1%), in (1.0%), of (0.3%), around (0.1%), ...</i>
<b>keep</b>	(69.8%)	<i>maintain (4.5%), hold (4.4%), be (4.2%), have (1.1%), make (1.0%), ...</i>
<b>his</b>	(86.2%)	<i>its (2.1%), statements (1.5%), what (1.0%), out (0.6%), the (0.6%), ...</i>
<b>statements</b>	(91.9%)	<i>testimony (1.5%), messages (0.7%), comments (0.6%), ...</i>
<b>vague</b>	(96.2%)	<i>v@@ (1.2%), in (0.6%), ambiguous (0.3%), ...</i>
<b>enough</b>	(98.9%)	<i>and (0.2%), ...</i>
<b>so</b>	(51.1%)	<i>, (44.3%), to (1.2%), in (0.6%), and (0.5%), just (0.2%), that (0.2%), ...</i>
<b>they</b>	(55.2%)	<i>that (35.3%), it (2.5%), can (1.6%), you (0.8%), we (0.4%), to (0.3%), ...</i>
<b>can</b>	(93.2%)	<i>may (2.7%), could (1.6%), are (0.8%), will (0.6%), might (0.5%), ...</i>
<b>be</b>	(98.4%)	<i>have (0.3%), interpret (0.2%), get (0.2%), ...</i>
<b>interpreted</b>	(99.1%)	<i>interpre@@ (0.1%), constru@@ (0.1%), ...</i>
<b>in</b>	(96.5%)	<i>on (0.9%), differently (0.5%), as (0.3%), to (0.2%), for (0.2%), by (0.1%), ...</i>
<b>different</b>	(41.5%)	<i>a (25.2%), various (22.7%), several (3.6%), ways (2.4%), some (1.7%), ...</i>
<b>ways</b>	(99.3%)	<i>way (0.2%), manner (0.2%), ...</i>
<b>.</b>	(99.2%)	<i>&lt;/s&gt; (0.2%), , (0.1%), ...</i>
<b>&lt;/s&gt;</b>	(100.0%)	



# ensembling

# Ensembling

- Train multiple models
- Say, by different random initializations

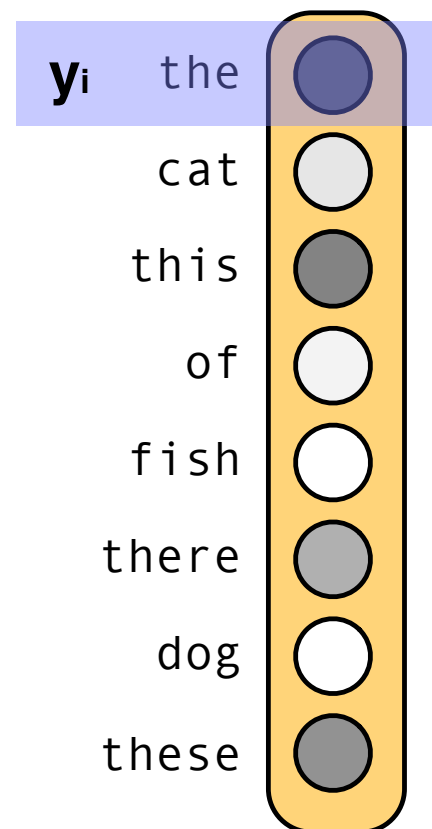
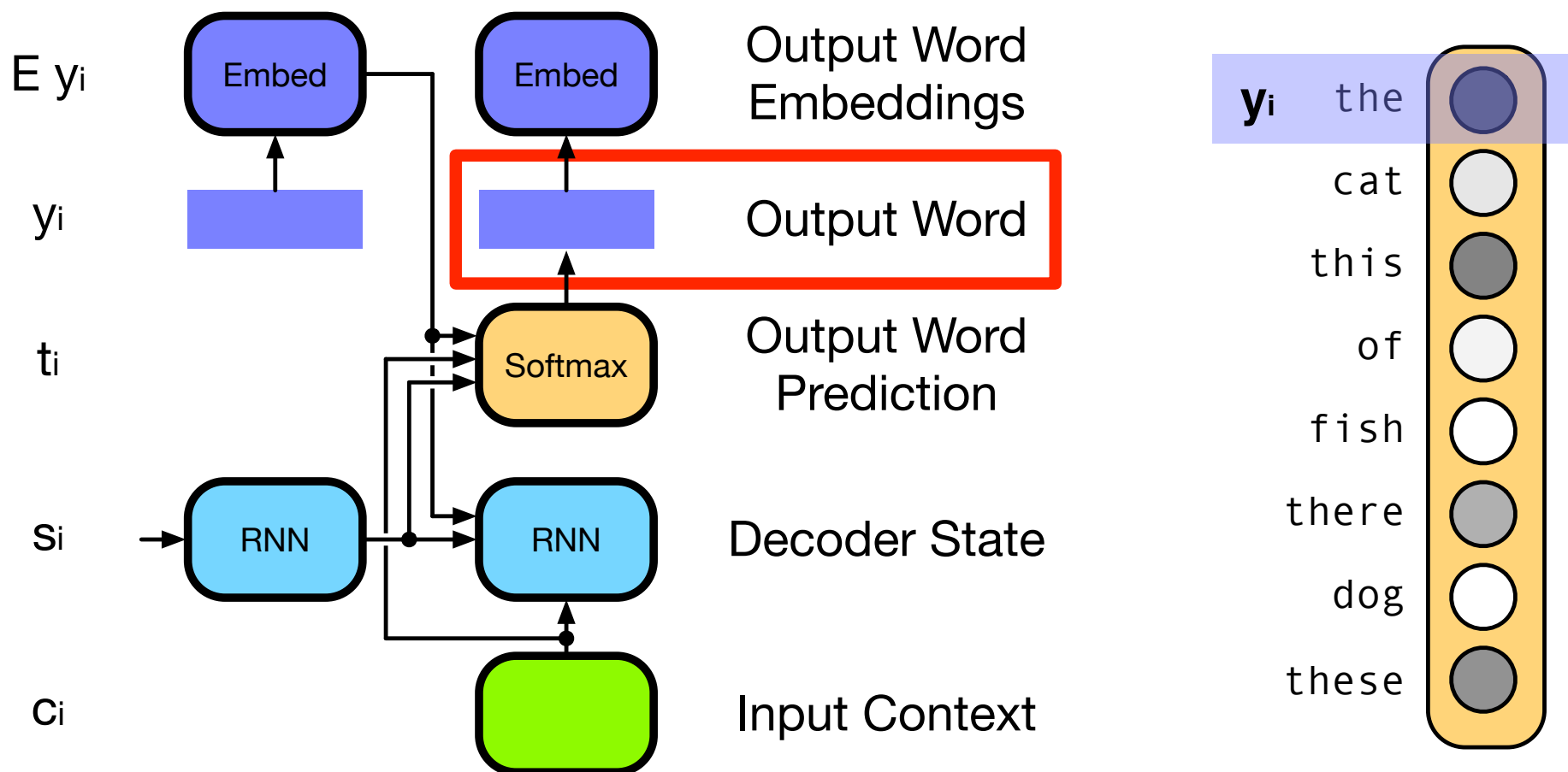


- Or, by using model dumps from earlier iterations

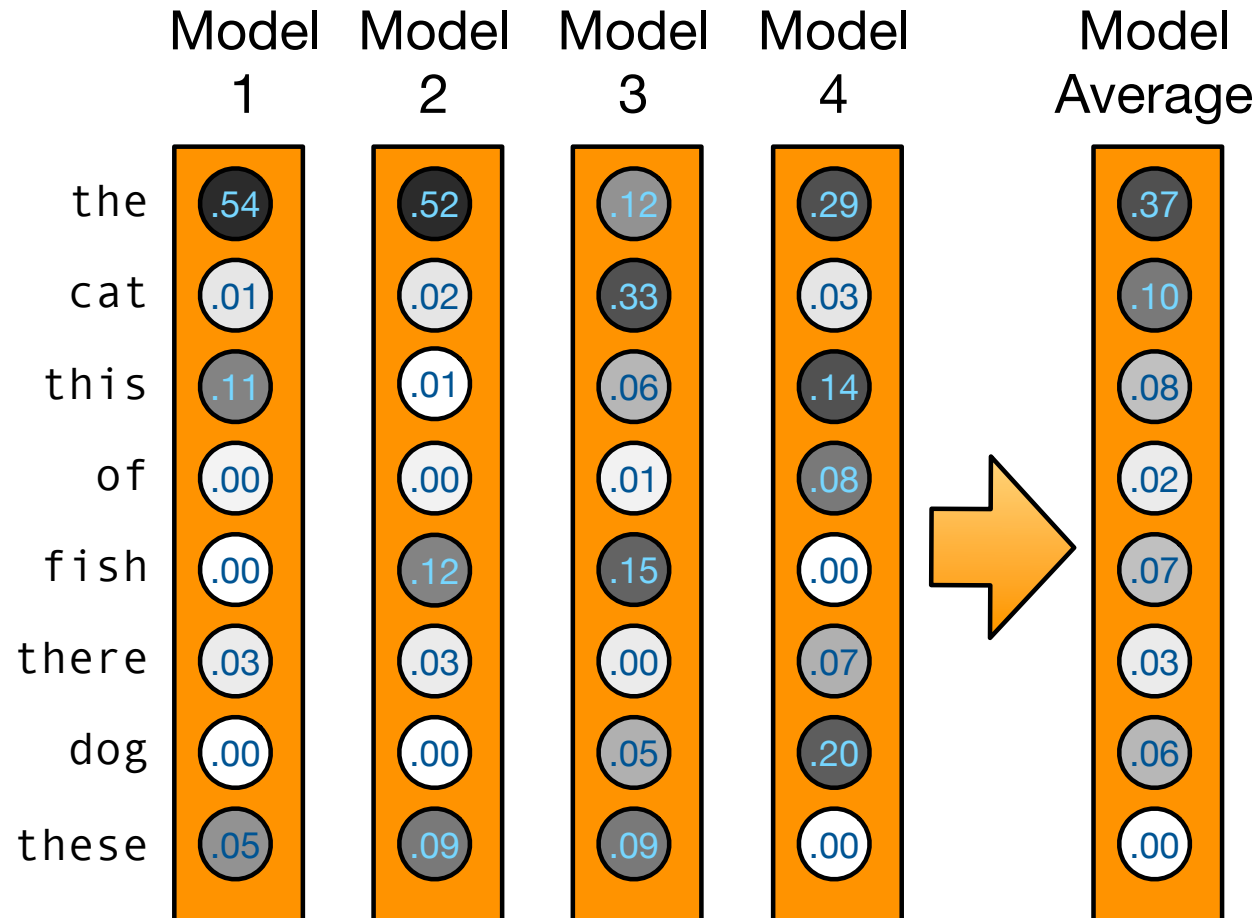


(most recent, or interim models with highest validation score)

# Decoding with Single Model



# Combine Predictions



# Ensembling

- Surprisingly reliable method in machine learning
- Long history, many variants:  
bagging, ensemble, model averaging, system combination, ...
- Works because errors are random, but correct decisions unique

# reranking

# Right-to-Left Inference

- Neural machine translation generates words right to left (L2R)

the → cat → is → in → the → bag → .

- But it could also generate them right to left (R2L)

the ← cat ← is ← in ← the ← bag ← .

**Obligatory notice:** Some languages (Arabic, Hebrew, ...) have writing systems that are right-to-left, so the use of "right-to-left" is not precise here.



# Right-to-Left Reranking

- Train both L2R and R2L model
- Score sentences with both
  - ⇒ use both left and right context during translation
- Only possible once full sentence produced → re-ranking
  1. generate n-best list with L2R model
  2. score candidates in n-best list with R2L model
  3. chose translation with best average score