

## Beam Search

Decoding strategies →

- LM's Place distribution  $P(y_i | y_1, \dots, y_{i-1})$
- Seq2Seq models place distribution  $P(y_i | x, y_1, \dots, y_{i-1})$
- Generation from both models looks similar; how do we do it?
  - Option 1: max  $y_i$   $P(y_i | y_1, \dots, y_{i-1})$  - take greedily best option
  - Option 2: Use beam search to find sequence w/ highest probability

Doesn't Work ✗ Option 3: Sample from model; draw  $y_i$  from that distribution

- Option 4: Nucleus Sampling

### • Beam Search

Assume vocab size  $|V|$  and sequence of length  $n$

$P(y_i)$	The	0.01	→	The dog	0.0004
	A	0.005	→	The cat	0.0003
	one	0.004	→	The fish	0.0002
	I	0.003	→		
	⋮	⋮			

End up with  $|V|^n$  total sequences  
↑  
too expensive

- Beam search: approximates exhaustive search w/ less compute
  - keep the top  $K$  hypotheses at each step

$K=3$	The	⊕	The dog	} End up with $K V $ total sequences
	A	⊕	The cat	
	one	⊕	A fish	
	one	⊕	one second	
	<del>I</del>		<del>X</del>	

### • Drawbacks of beam search:

- Degeneration: If a fragment is repeated 2-3x, it has a high probability to keep repeating. Model fails because it's locally normalized

Runtime:  $K \times n$  transformer calls

$K \times |V| \times n$  hypotheses

### • Drawbacks of Sampling: Long Tail

$P(y | \text{they live in a desert near})$

0.01	roads	} Good options, account for 90% of total probability mass, 90% chance of good prediction
0.01	towns	
0.01	people	
0.005	struts	
...		
0.0005	town	- Long tail w/ 10% of mass

solution →

Nucleus Sampling

\* On average, every 10 words, we get something from 10% tail of distribution. We don't want this!