

DECODING TECHNIQUES

SUBTITLE

Ramaseshan Ramachandran

PRE-TRAINING

- ✦ Process of learning linguistic patterns & world knowledge from massive text datasets
- ✦ Trains models to predict text sequences, building foundational language understanding

LLMS

- ✦ Output model from Pre-training - severe billion parameters
- ✦ Develop emergent capabilities through scale:
- ✦ Contextual word representations
- ✦ Cross-domain knowledge retention
- ✦ Pattern recognition across languages

ADVANTAGES OF LLMS

- ✦ State-of-the-art performance on NLP benchmarks:
 - ✦ Text generation (most transformative)
 - ✦ Semantic understanding
 - ✦ Few-shot learning
 - ✦ Particularly effective for generative tasks
 - ✦ Summarization
 - ✦ Machine Translation
 - ✦ Question Answering
 - ✦ Chatbot Interactions

DECODING TECHNIQUES

- ✦ Selecting next token from probability distribution
- ✦ Key components:
 - ✦ Context window (prior generated text)
 - ✦ Vocabulary probability scores
 - ✦ Decoding strategy algorithm
- ✦ Repeatedly choosing the next word conditioned on the previous choices - autoregressive/causal generation

RANDOM SAMPLING

- ✦ Generates sensible, high-probability words but also includes odd, low-probability words, resulting in weird sentences
- ✦ Will it effectively generate sentences with adequate and fluent structure?
- ✦ We look for quality and diversity in the generated text
- ✦ We want techniques that emphasize the most probable words

GREEDY APPROACH

- Model is computed using conditional probabilities. We want to generate a sentence $w_1, w_2, w_3, \dots, w_n$ using

- $$\hat{w}_i = \arg \max_{w \in V} P(w_i | w_{<i})$$

- The approach makes a locally optimal choices - Highest probability token is chosen
- Generate words that are likely in the context and less likely to generate equivalent words that are unlikely

- Generates sentences that are
 - More accurate
 - More coherent
 - More factual
 - Boring and more repetitive
- What happens if we choose the next word from the the middle of the distribution
 - May be more creative and diverse
 - Less factual and incoherent and no adequate

TEMPERATURE

- ✦ LLM Temperature Impact: Significantly affects text coherence
 - ✦ Controls token selection randomness
- ✦ Quality of decoding: Impacts the quality of the output generated by the LLM - Greedy → Balanced choices

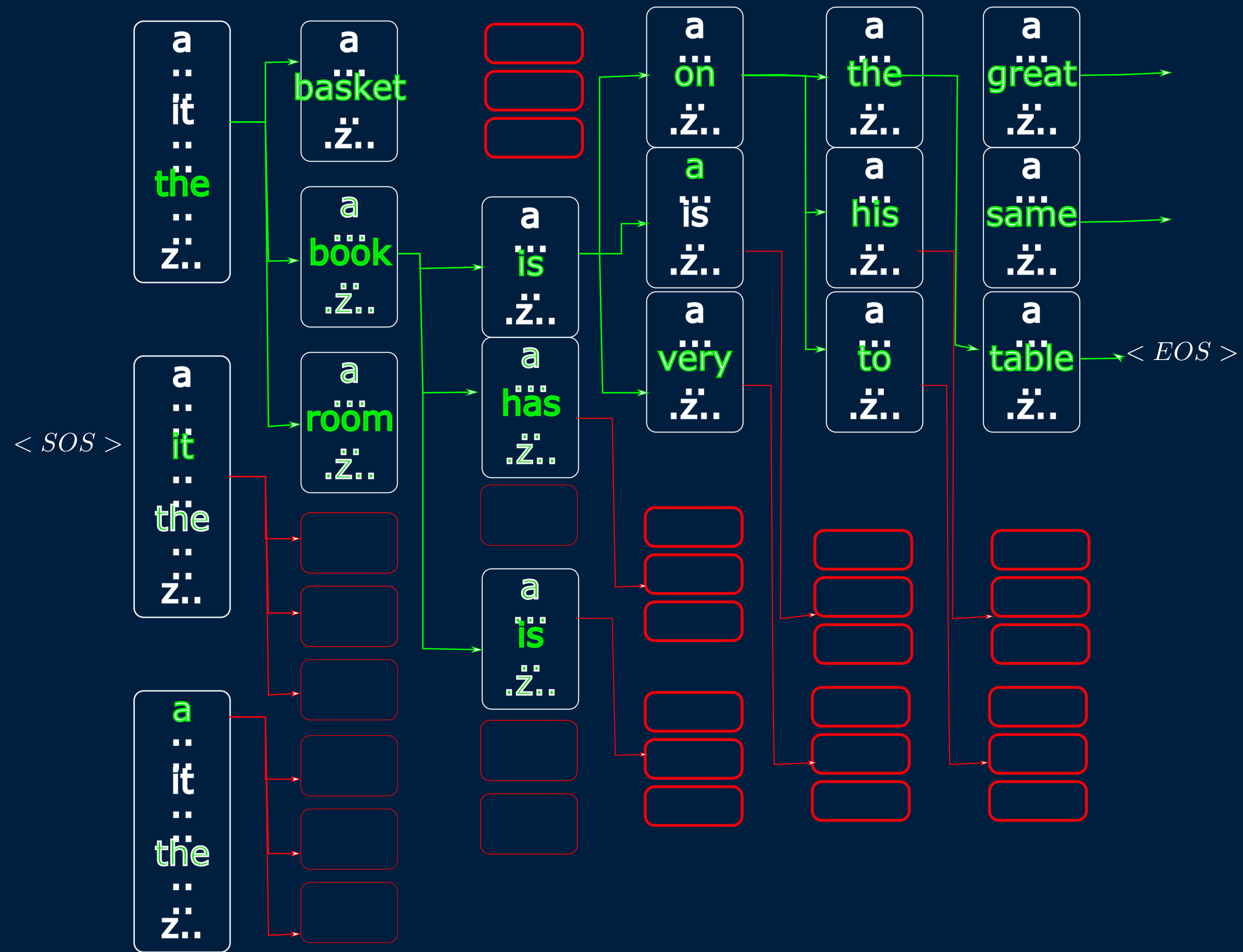
$$\star P(w_t) = \frac{\exp(z_i/T)}{\sum \exp(z_j/T)}$$

- ✦ $z = \mathbf{h} \cdot \mathbf{W}_{\text{vocab}} + \mathbf{b}$, represents raw *logit* score before applying any activation function

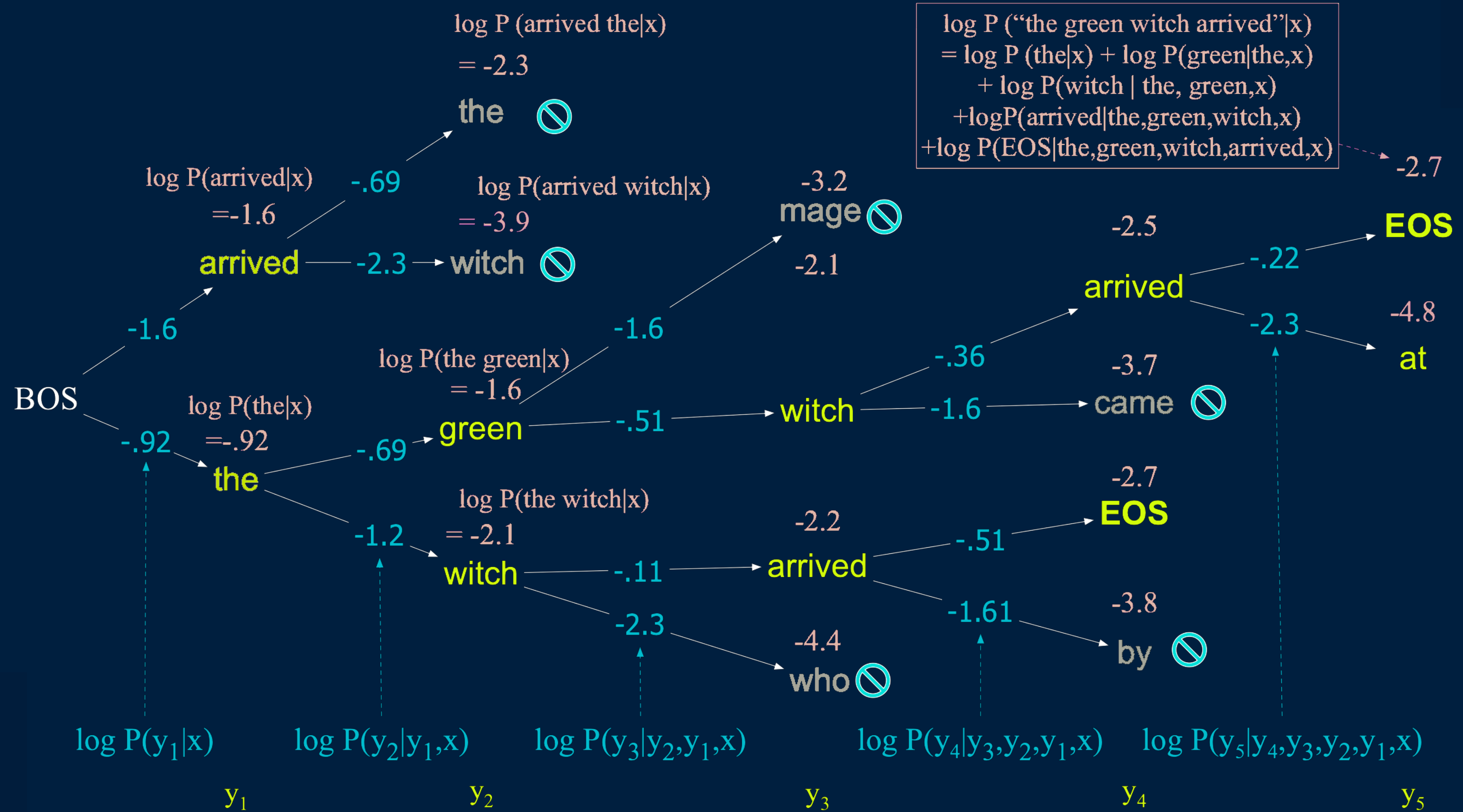
BEAM SEARCH I

- ✦ Selects a few candidate hypothesis from $|V|$. It reduces memory requirement by using only a $M < |V|$ candidates using a score.
 - ✦ Maintain M candidates/hypothesis at each time step
 - ✦ $C_t = (x_1^1, \dots, x_t^1) \dots (x_1^M, \dots, x_t^M)$
 - ✦ Compute C_{t+1} by expanding C_t and keeping the best M candidates
 - ✦ $\tilde{C} = \bigcup_{i=1}^M C_{t-1}^i$
- ✦ Typical Beam width = 5-10

BEAM SEARCH II



BEAM EXAMPLE ($M = 2$)



TOP-K SAMPLING

- ✦ Truncate the distribution to the *top-k* most likely words.
- ✦ Renormalized to produce a probability distribution
- ✦ A word is randomly sampled from within the *top-k* words according to their renormalized probabilities
- ✦ When $k = 1$, *top-k* sampling is identical to greedy decoding
- ✦ Setting k to a larger number
 - ✦ More diverse but still high-quality text
 - ✦ Impact on fluency - Low risk
- ✦ Selecting to the middle-probability words
 - ✦ More creative and more diverse
 - ✦ Impact on fluency - High risk

Use case	Top-k	Rationale
Technical writing	5-10	High-confidence tokens
Creative writing	20-50	Controlled generation for creativity
Conversational AI	10-30	Safety and engagement

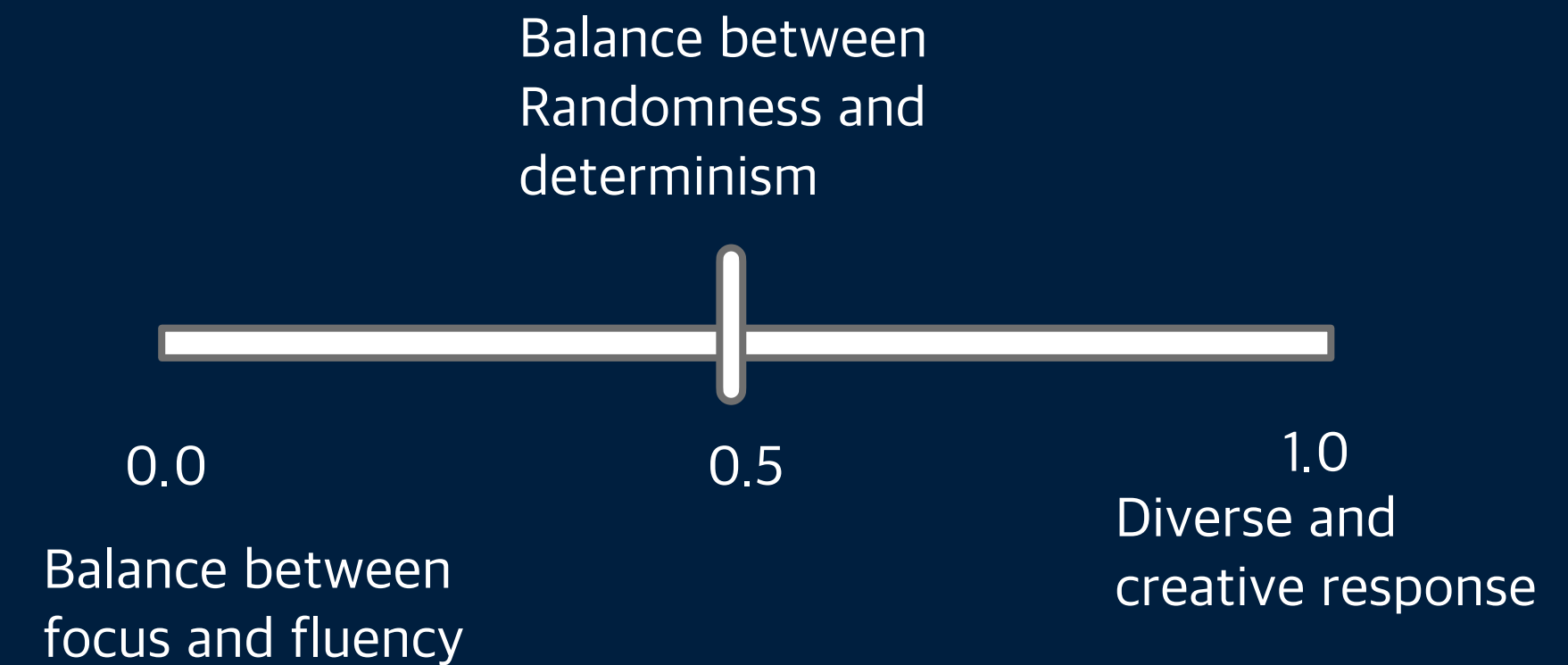
N. N. Minh, A. Baker, C. Neo, A. G. Roush, A. Kirsch, and R. Schwartz-Ziv. Turning up the heat:Min-p sampling for creative and coherent LLM outputs. In The Thirteenth International Conference on Learning Representations, 2025.

TOP-P SAMPLING

- ✦ Keeps the *top-p* percent of the probability mass
 - ✦ Tokens selected = $\arg \min_k \sum_{i=1}^k P(w_i) \leq p$
 - ✦ Selects the smallest set of tokens whose cumulative probability exceeds threshold p
 - ✦ Removes very unlikely words
- ✦ Measures probability rather than the number of words
- ✦ Balances creativity and coherence by dynamically adjusting the candidate pool
- ✦ Rescales probability of the selected tokens so that

$$\sum_i p_i = 1$$

ChatGpt uses top-p sampling



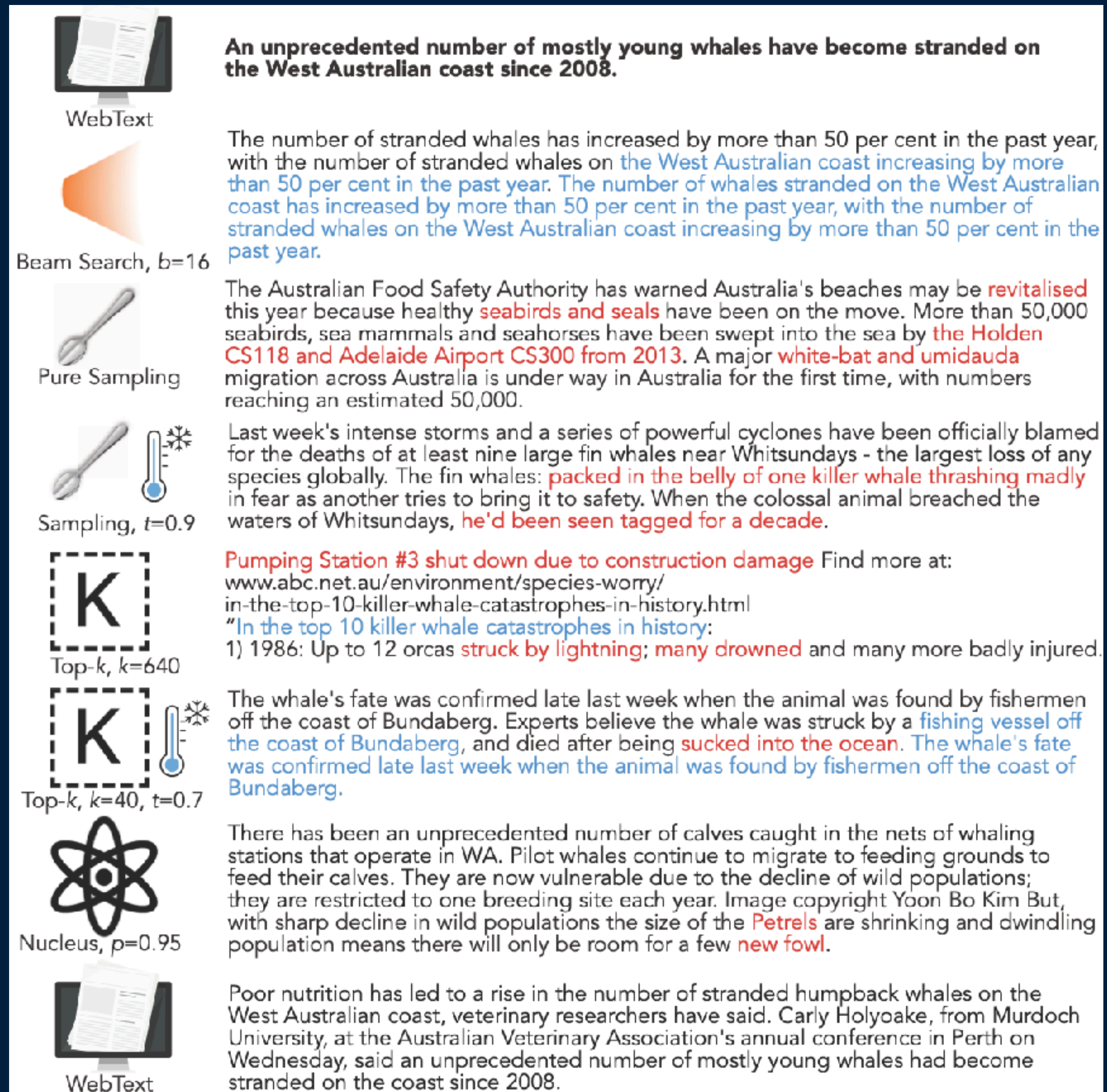


Figure 3: Example generations continuing an initial sentence. Maximization and top- k truncation methods lead to copious repetition (highlighted in blue), while sampling with and without temperature tends to lead to incoherence (highlighted in red). Nucleus Sampling largely avoids both issues.

MIN-P SAMPLING

- ✦ Adjusts cutoff threshold based on model confidence in real-time
 - ✦ $\text{Threshold}_t = \max(P(x_t | x_{1:t-1})) \times \text{min_p}$ Where:
 - ✦ $\max(P(x_t | x_{1:t-1}))$ = highest token probability at step
 - ✦ min_p = user-defined ratio (e.g., 0.05-0.2)
- ✦ At each generation step:
 - ✦ Compute token probabilities $P(x_t | x_{1:t-1})$ over vocabulary V
 - ✦ Identify top probability p_{\max}
 - ✦ Calculate adaptive threshold:
 - ✦ $\text{Threshold} = p_{\max} \times \text{min_p}$. Retain tokens where $p_i \geq \text{Threshold}$
 - ✦ Sample from filtered distribution

TEMPERATURE BASED SAMPLING

♦ T=0.2: [██████████] Top token - *dominates*

♦ T=1.0: [██████████] Balanced - *exploration*

♦ T=2.0: [██████████] High-risk - *diversity*

♦ Use Temperature When:

- ♦ Creativity response is required (e.g., brainstorming, storytelling)
- ♦ Balancing exploration/exploitation (e.g., chatbots)

♦ Avoid Temperature When:

- ♦ Maximum determinism is required (e.g., legal contracts)
- ♦ Using pure greedy/beam search

Decoding approach	Temperature used
Greedy	No
Beam search	No
Temperature Sampling	Yes
Top-K	Yes
Top-P	Yes
Min-P	Yes