

# CS 6501: Constrained-Aware Generative AI

## Lecture Notes 4

### Autoregressive Transformers and Decoding

*Generation as search under probabilistic and symbolic constraints*

Prof. Ferdinando Fioretto

Department of Computer Science, University of Virginia

Thursday, January 22, 2026

#### Abstract

Autoregressive Transformers are the dominant foundation for modern language models because they couple a tractable likelihood with a powerful sequence architecture. In constrained-aware generation, this matters for two reasons. First, the factorization yields a canonical probabilistic objective, so “validity” can be expressed as a distributional constraint that modifies a base model. Second, inference at generation time is an explicit *search* over token sequences, so constraints can be injected directly into decoding by pruning invalid continuations and reweighting feasible candidates. These notes review the autoregressive factorization, the Transformer architecture, and standard decoding methods, then introduce constrained decoding via grammars and finite-state constraints as the first concrete instance of “generation as search subject to constraints”.

## 1 From constrained generation to constrained decoding

Lecture 1 introduced a recurring template for constrained-aware generation: starting from a base model  $p_{\theta}(\mathbf{x} \mid \mathbf{c})$  that encodes plausibility, and combining it with soft penalties and hard feasibility requirements to define a constrained target distribution  $\pi(\mathbf{x} \mid \mathbf{c})$ . Autoregressive language models specialize this view to sequences  $\mathbf{y} = (y_1, \dots, y_T)$  over a vocabulary  $\mathcal{V}$ , conditioned on a prompt or context  $\mathbf{x}$  (which may itself be a token sequence). The central algorithmic fact for this lecture is that, for autoregressive models, *generation is an incremental search over prefixes*. This opens a direct insertion point for constraints that act on partial sequences.

#### Constrained decoding as constrained inference

Given a prompt  $\mathbf{x}$ , an autoregressive model defines  $p_{\theta}(\mathbf{y} \mid \mathbf{x})$ . A hard constraint set  $\mathcal{C}(\mathbf{x}) \subseteq \mathcal{V}^T$  induces the constrained distribution

$$\pi(\mathbf{y} \mid \mathbf{x}) \propto p_{\theta}(\mathbf{y} \mid \mathbf{x}) \mathbf{1}\{\mathbf{y} \in \mathcal{C}(\mathbf{x})\}, \quad (1)$$

and a soft constraint  $\phi(\mathbf{y}, \mathbf{x})$  yields the energy-shaped target

$$\pi(\mathbf{y} \mid \mathbf{x}) \propto p_{\theta}(\mathbf{y} \mid \mathbf{x}) \exp(-\lambda \phi(\mathbf{y}, \mathbf{x})). \quad (2)$$

Decoding is any procedure that approximates either sampling  $\mathbf{y} \sim \pi(\cdot \mid \mathbf{x})$  or optimization  $\hat{\mathbf{y}} \in \arg \max_{\mathbf{y}} \log \pi(\mathbf{y} \mid \mathbf{x})$ .

The remainder of the lecture makes the above concrete by (i) recalling how  $p_{\theta}(\mathbf{y} \mid \mathbf{x})$  is defined and trained for Transformers, (ii) formalizing standard decoding rules, and (iii) showing how grammars and finite-state constraints can be enforced at the level of prefixes.

## 2 Autoregressive factorization and maximum likelihood

Let  $\mathbf{y} = (y_1, \dots, y_T)$  be a token sequence. The autoregressive assumption models the joint distribution as

$$p_{\theta}(\mathbf{y}) = \prod_{t=1}^T p_{\theta}(y_t | y_{<t}), \quad \text{where} \quad y_{<t} := (y_1, \dots, y_{t-1}). \quad (3)$$

In conditional form, given a prompt  $\mathbf{x}$ , the factorization becomes

$$p_{\theta}(\mathbf{y} | \mathbf{x}) = \prod_{t=1}^T p_{\theta}(y_t | y_{<t}, \mathbf{x}). \quad (4)$$

Training by maximum likelihood on a dataset  $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n$  corresponds to minimizing the negative log-likelihood, equivalently the token-level cross-entropy:

$$\min_{\theta} \mathcal{L}(\theta) := - \sum_{i=1}^n \sum_{t=1}^{T_i} \log p_{\theta}\left(y_t^{(i)} | y_{<t}^{(i)}, \mathbf{x}^{(i)}\right). \quad (5)$$

Here  $T_i$  counts predicted tokens up to (and including) EOS but excludes padding; in a decoder-only setting the prompt  $\mathbf{x}$  is typically concatenated ahead of  $\mathbf{y}$  so that  $y_{<t}$  includes the prompt tokens. This objective is optimized with teacher forcing: during training, the model conditions on the *true* prefix  $y_{<t}^{(i)}$ . At test time, the model conditions on its own generated prefix, which yields exposure bias and compounds errors [Bengio et al., 2015, Ranzato et al., 2016]. From the viewpoint of constrained-aware generation, this gap is important: if constraints are not enforced during decoding, small local errors can push the model into regions where global validity is hard to recover.

A standard evaluation summary of (5) is perplexity. For a test set, define the average per-token negative log-likelihood  $\ell$  as (5) divided by the total number of tokens. The perplexity is then  $\text{PPL} := \exp(\ell)$ . Perplexity is useful as a calibration measure of probabilistic fit, but it is not a validity guarantee: a low-perplexity continuation may still violate grammar constraints, lexical requirements, or policy constraints [Theis et al., 2016, Holtzman et al., 2020].

## 3 Transformers for autoregressive modeling

The autoregressive factorization (4) reduces conditional generation to learning a family of next-token conditionals  $p_{\theta}(y_t | y_{<t}, \mathbf{x})$  that can be evaluated efficiently for all prefixes in parallel and updated as the prefix grows. At this point the modeling question is: which parameterization should we choose for these conditional distributions? In principle, many sequence models can serve this role, including  $n$ -gram, recurrent networks, and state-space models. However, in modern large-scale language modeling, the *Transformer* architecture has emerged as the dominant parameterization of  $p_{\theta}$  because it supports stable optimization at scale, exposes a highly expressive attention-based representation of the context  $(y_{<t}, \mathbf{x})$ , and can be implemented with hardware-efficient batched matrix operations [Vaswani et al., 2017]. We therefore focus on autoregressive Transformers as the canonical instantiation of (4), and then study decoding as the induced search problem for high-probability sequences, including the first constrained variants.

**Token and position representations.** Let an input sequence of length  $L$  be represented by token ids  $(u_1, \dots, u_L)$  with  $u_i \in \mathcal{V}$ , where  $\mathcal{V}$  is a finite vocabulary and  $|\mathcal{V}|$  denotes its size. A Transformer first maps each token id to a dense vector in  $\mathbb{R}^d$  through an embedding table

$E \in \mathbb{R}^{|\mathcal{V}| \times d}$ , so that  $e_i := E[u_i] \in \mathbb{R}^d$ . Because attention is permutation-invariant in its basic form, the model must also encode the position index  $i \in \{1, \dots, L\}$ . This is done by adding a positional encoding  $p_i \in \mathbb{R}^d$ , either learned (via a separate table  $P \in \mathbb{R}^{L_{\max} \times d}$ ) or fixed (e.g., sinusoidal). The initial hidden state at position  $i$  is then

$$h_i^{(0)} := e_i + p_i \in \mathbb{R}^d.$$

Stacking the vectors row-wise yields a matrix of hidden states  $H^{(0)} \in \mathbb{R}^{L \times d}$  whose  $i$ -th row is  $h_i^{(0)\top}$ .

**Self-attention with a causal mask.** The central operation in a Transformer layer is to update each position's representation by taking a content-dependent weighted average of representations at other positions. Fix a layer index  $\ell \in \{1, \dots, L_{\text{layers}}\}$  and suppose the previous layer provides hidden states  $H^{(\ell-1)} \in \mathbb{R}^{L \times d}$ , where the  $i$ -th row  $h_i^{(\ell-1)\top}$  summarizes the prefix information available to position  $i$  at depth  $\ell - 1$ .

Self-attention is parameterized via three linear maps applied to every position: the *queries*, *keys*, and *values*. Concretely, choose a head dimension  $d_k$  and define learned matrices  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d_k}$ . The projected matrices are

$$Q := H^{(\ell-1)} W_Q \in \mathbb{R}^{L \times d_k}, \quad K := H^{(\ell-1)} W_K \in \mathbb{R}^{L \times d_k}, \quad V := H^{(\ell-1)} W_V \in \mathbb{R}^{L \times d_k}.$$

Write  $q_i^\top$  for the  $i$ -th row of  $Q$  and  $k_j^\top$  for the  $j$ -th row of  $K$ . The scalar compatibility score between position  $i$  (as a query) and position  $j$  (as a key) is the dot product  $\langle q_i, k_j \rangle = q_i^\top k_j$ . Collecting all pairwise scores yields the matrix  $S := QK^\top \in \mathbb{R}^{L \times L}$ , where  $S_{ij} = q_i^\top k_j$ . The scaled dot-product attention operator is

$$\text{Attn}(Q, K, V) = AV \quad \text{with} \quad A := \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}} + M\right) \in \mathbb{R}^{L \times L}, \quad (6)$$

where  $\sqrt{d_k}$  is a normalization that stabilizes magnitudes of dot products as  $d_k$  grows, and the softmax is applied row-wise, i.e., for each  $i$ ,

$$A_{ij} = \frac{\exp\left(\frac{S_{ij}}{\sqrt{d_k}} + M_{ij}\right)}{\sum_{r=1}^L \exp\left(\frac{S_{ir}}{\sqrt{d_k}} + M_{ir}\right)}.$$

Thus, the updated representation at position  $i$  is a convex combination of value vectors: if  $v_j^\top$  is the  $j$ -th row of  $V$ , then the  $i$ -th row of  $\text{Attn}(Q, K, V)$  equals  $\sum_{j=1}^L A_{ij} v_j$ .

For autoregressive language modeling, attention must be *causal*: position  $i$  is allowed to depend only on positions  $j \leq i$ . This is enforced via an additive mask matrix  $M \in \mathbb{R}^{L \times L}$  defined by

$$M_{ij} := \begin{cases} 0, & j \leq i, \\ -\infty, & j > i, \end{cases}$$

so that after adding  $M$  and applying the softmax, all weights  $A_{ij}$  with  $j > i$  become 0. In practice,  $-\infty$  is implemented as a large negative constant, and padding positions can be masked by setting the corresponding entries in  $M$  to  $-\infty$  as well, ensuring that attention ignores padded tokens.

*Multi-head attention* increases expressivity by running several attention operators in parallel with different learned projections. With  $H$  heads, one typically chooses  $d_k = d/H$  so the total width remains  $d$ . For head  $h \in \{1, \dots, H\}$ , define  $W_Q^{(h)}, W_K^{(h)}, W_V^{(h)} \in \mathbb{R}^{d \times d_k}$ , compute  $\text{Attn}^{(h)} \in \mathbb{R}^{L \times d_k}$  via (6), and then concatenate head outputs along the feature dimension:

$$\text{MHA}(H^{(\ell-1)}) := [\text{Attn}^{(1)} \| \dots \| \text{Attn}^{(H)}] W_O,$$

where  $\parallel$  denotes concatenation, the bracketed term lies in  $\mathbb{R}^{L \times d}$ , and  $W_O \in \mathbb{R}^{d \times d}$  is an output projection. Intuitively, different heads can specialize to different relations, such as local syntactic dependencies, long-range coreference, or delimiter matching that is relevant for constrained generation.

A standard decoder block then combines attention with a position-wise feedforward network (FFN) and stabilization layers. A common choice is the two-layer FFN

$$\text{FFN}(x) := W_2 \sigma(W_1 x + b_1) + b_2,$$

applied independently to each position, where  $W_1 \in \mathbb{R}^{d \times d_{\text{ff}}}$ ,  $W_2 \in \mathbb{R}^{d_{\text{ff}} \times d}$ ,  $d_{\text{ff}}$  is an expansion dimension, and  $\sigma$  is a nonlinearity (e.g., GELU). Residual connections and normalization (LayerNorm or RMSNorm) are applied around both the attention and FFN sublayers; in the common pre-norm form,

$$\tilde{H} = H^{(\ell-1)} + \text{MHA}(\text{Norm}(H^{(\ell-1)})), \quad H^{(\ell)} = \tilde{H} + \text{FFN}(\text{Norm}(\tilde{H})).$$

These design choices are not cosmetic: they enable training deep stacks and are central to why Transformers scale effectively [Vaswani et al., 2017].

**From hidden states to next-token probabilities.** In a decoder-only language model, the final layer outputs a sequence of hidden states  $H^{(L_{\text{layers}})} \in \mathbb{R}^{L \times d}$ . At a time step  $t$  (corresponding to position  $t$  in the concatenated prompt-plus-generated sequence), the model maps the hidden vector  $h_t := h_t^{(L_{\text{layers}})} \in \mathbb{R}^d$  to logits over the vocabulary via an output matrix  $W_U \in \mathbb{R}^{d \times |\mathcal{V}|}$  and (optionally) a bias  $b_U \in \mathbb{R}^{|\mathcal{V}|}$ :

$$s_t := h_t^\top W_U + b_U \in \mathbb{R}^{|\mathcal{V}|}.$$

The component  $s_{t,v}$  is the unnormalized score assigned to token  $v \in \mathcal{V}$  at step  $t$ . The conditional distribution is then obtained by the softmax,

$$p_\theta(y_t = v \mid y_{<t}, \mathbf{x}) = \frac{\exp(s_{t,v})}{\sum_{v' \in \mathcal{V}} \exp(s_{t,v'})}. \quad (7)$$

Often  $W_U$  is tied to the embedding table  $E$  (weight tying), which reduces parameters and can improve generalization. Under scaling and sufficient data, this parameterization supports few-shot prompting and in-context learning behaviors [Brown et al., 2020, Wei et al., 2022].

**Encoder-only and encoder-decoder variants.** While this course emphasizes autoregressive decoding, it is useful to contrast other Transformer families. BERT-style models are encoder-only: they process the entire input bidirectionally and are trained with masked language modeling, a denoising objective that predicts masked tokens from both left and right context [Devlin et al., 2019]. Sequence-to-sequence models adopt an encoder-decoder structure: an encoder maps an input  $\mathbf{x}$  to contextual states, and a decoder generates  $\mathbf{y}$  autoregressively while attending to the encoder states via cross-attention. Historically, many constrained generation problems in translation and summarization were studied in this setting because the separation between source conditioning and target decoding makes it natural to incorporate lexical and structural constraints during decoding.

## 4 Pretraining, scaling, and compute-data tradeoffs

Modern language models are typically trained in two stages: unsupervised pretraining on large corpora, then task adaptation (supervised fine-tuning, instruction tuning, or preference optimization). Several empirical regularities matter for constrained-aware generation.

Kaplan et al. proposed scaling laws that relate loss to model size, dataset size, and compute under regimes where optimization is stable [Kaplan et al., 2020]. Hoffmann et al. argued that, for a fixed compute budget, many models were under-trained on data, and derived compute-optimal tradeoffs that informed later LLM training recipes [Hoffmann et al., 2022]. These scaling perspectives are directly relevant to constraint-aware generation because constraint satisfaction can be viewed as requiring *effective capacity* on rare events: hard constraints often correspond to low-probability regions under the base model, so feasibility under constraints may demand either additional model capacity, additional data that covers constrained regions, or explicit inference-time constraint mechanisms.

Open technical reports such as GPT-2 [Radford et al., 2019] and LLaMA [Touvron et al., 2023] provide useful engineering context, including architectural choices, tokenization, and training setups. However, for this course the key point is conceptual: even very large base models remain stochastic predictors, so correctness and validity often require explicit constraint enforcement during decoding, reranking, or post-processing.

## 5 Decoding: sampling and search in autoregressive models

Given a prompt  $\mathbf{x}$  and a model  $p_{\theta}(\mathbf{y} \mid \mathbf{x})$ , generation is an iterative procedure that chooses tokens  $y_t$  sequentially. Formally, at each step  $t$  we have a distribution over next tokens as in (7), and we select a token by either (i) sampling, or (ii) approximate maximization of some global objective.

### 5.1 An inference view of decoding

A common objective is the maximum a posteriori (MAP) continuation

$$\hat{\mathbf{y}} \in \arg \max_{\mathbf{y} \in \mathcal{V}^T} [\log p_{\theta}(\mathbf{y} \mid \mathbf{x}) + \alpha \cdot \text{len}(\mathbf{y})], \quad (8)$$

where the length term (or an explicit EOS penalty) handles variable-length decoding. Alternatively, one may aim to sample from  $p_{\theta}(\cdot \mid \mathbf{x})$  to preserve diversity. For open-ended generation, naive ancestral sampling can lead to repetitive and degenerate text [Holtzman et al., 2020], motivating a family of heuristic decoding rules that shape the next-token distribution.

### 5.2 Greedy decoding and temperature sampling

Greedy decoding chooses

$$y_t \in \arg \max_{v \in \mathcal{V}} \log p_{\theta}(v \mid y_{<t}, \mathbf{x}). \quad (9)$$

This is fast but often yields generic continuations.

Temperature sampling rescales logits as  $s_{t,v}/\tau$  for  $\tau > 0$ , defining

$$p_{\theta}^{(\tau)}(v \mid y_{<t}, \mathbf{x}) := \frac{\exp(s_{t,v}/\tau)}{\sum_{v' \in \mathcal{V}} \exp(s_{t,v'}/\tau)}. \quad (10)$$

As  $\tau \downarrow 0$ , the distribution concentrates near the greedy choice; as  $\tau \uparrow \infty$ , it approaches uniform over the vocabulary. Temperature is therefore a knob that trades off diversity and local likelihood.

### 5.3 Top- $k$ and nucleus (top- $p$ ) sampling

Top- $k$  sampling restricts attention to the  $k$  most likely tokens at step  $t$ :

$$\mathcal{V}_t^{(k)} := \text{top-}k \text{ tokens under } p_{\theta}(\cdot \mid y_{<t}, \mathbf{x}), \quad p(v) \leftarrow \frac{p(v) \mathbf{1}\{v \in \mathcal{V}_t^{(k)}\}}{\sum_{v' \in \mathcal{V}_t^{(k)}} p(v')}. \quad (11)$$

**Algorithm 1** Nucleus (top- $p$ ) sampling

---

```

1: Input: prompt  $\mathbf{x}$ , model  $p_{\theta}$ , threshold  $p \in (0, 1)$ , temperature  $\tau > 0$ 
2: Initialize prefix  $\mathbf{y}_{<1} \leftarrow \emptyset$ 
3: for  $t = 1, 2, \dots$  until EOS or length limit do
4:   Compute logits  $s_t$  from the Transformer given  $(\mathbf{x}, \mathbf{y}_{<t})$ 
5:   Form probabilities  $q(v) \propto \exp(s_{t,v}/\tau)$  over  $v \in \mathcal{V}$ 
6:   Sort tokens by  $q(v)$  and let  $\mathcal{V}_t^{(p)}$  be the smallest set with  $\sum_{v \in \mathcal{V}_t^{(p)}} q(v) \geq p$ 
7:   Renormalize  $\tilde{q}(v) \propto q(v) \mathbf{1}\{v \in \mathcal{V}_t^{(p)}\}$  and sample  $y_t \sim \tilde{q}(\cdot)$ 
8:   Append  $y_t$  to the prefix  $\mathbf{y}_{<t+1} \leftarrow (\mathbf{y}_{<t}, y_t)$ 
9: end for
10: Return: generated sequence  $\mathbf{y}$ 
```

---

Nucleus sampling (top- $p$ ) chooses the smallest set  $\mathcal{V}_t^{(p)}$  whose cumulative probability mass is at least  $p \in (0, 1)$ , then renormalizes similarly [Holtzman et al., 2020]. Unlike top- $k$ , top- $p$  adapts the candidate set size to the entropy of the distribution: in sharp distributions it behaves like greedy, in flat distributions it keeps more candidates.

the candidate set size to the entropy of the distribution: in sharp distributions it behaves like greedy, in flat distributions it keeps more candidates.

These heuristics can be interpreted as approximate constrained sampling with a *stepwise* feasibility set. However, because the truncation is local, it does not enforce global properties such as satisfying a grammar or including specific required substrings. This motivates explicit constrained decoding methods.

## 5.4 Beam search as approximate MAP

Beam search approximates (8) by maintaining a set of  $B$  partial hypotheses at each step and extending each hypothesis by likely next tokens. Let  $b$  index hypotheses and let  $\ell(b)$  denote a score, often the accumulated log-probability (possibly length-normalized). Beam search proceeds by expanding each hypothesis, then keeping the  $B$  best scoring continuations.

Beam search is a generic approximate search method, not an exact inference procedure. In open-ended generation, beam search can amplify model biases toward high-probability generic text. In constrained settings, beam search is attractive because constraints can be injected by restricting allowable expansions, so the beam never enters infeasible regions.

## 5.5 Reranking and decoding as optimization

A practical pattern is two-stage generation: first produce a candidate set  $\{\mathbf{y}^{(m)}\}_{m=1}^M$  using sampling or beam search, then choose the best candidate under an auxiliary scoring function:

$$\hat{\mathbf{y}} \in \arg \max_{m \in [M]} \left[ \log p_{\theta}(\mathbf{y}^{(m)} \mid \mathbf{x}) - \lambda \phi(\mathbf{y}^{(m)}, \mathbf{x}) \right]. \quad (12)$$

Here  $\phi$  may be a constraint violation score, a verifier output, a toxicity detector, or an external simulator. This is a first instance of “verifier-in-the-loop” constrained generation.

Contrastive decoding [Li et al., 2022] can be interpreted in this family: it modifies the token selection rule by preferring tokens that are simultaneously likely under an expert model and unlikely under a weaker or noisier model, biasing generation away from generic high-frequency patterns. The optimization perspective is useful for this course because it connects decoding rules to constraint energies, which later reappear in guidance, projection, proximal updates, and differentiable layers.

**Algorithm 2** FSA-constrained decoding (generic skeleton)

---

```

1: Input: prompt  $\mathbf{x}$ , model  $p_{\theta}$ , automaton  $A = (\mathcal{S}, \mathcal{V}, \delta, s_0, \mathcal{F})$ , decoding rule Select
2: Initialize prefix  $\mathbf{y}_{<1} \leftarrow \emptyset$  and automaton state  $s_1 \leftarrow s_0$ 
3: for  $t = 1, 2, \dots$  until EOS or length limit do
4:   Compute next-token distribution  $p_{\theta}(\cdot | \mathbf{y}_{<t}, \mathbf{x})$ 
5:   Compute allowed set  $\mathcal{A}(s_t) = \{v : \delta(s_t, v) \neq \emptyset\}$ 
6:   Mask and renormalize to obtain  $\pi(\cdot | \mathbf{y}_{<t}, \mathbf{x}, s_t)$  as in (15)
7:   Choose  $y_t \leftarrow \text{Select}(\pi(\cdot | \mathbf{y}_{<t}, \mathbf{x}, s_t))$  (sampling, greedy, or beam expansion)
8:   Update automaton state  $s_{t+1} \in \delta(s_t, y_t)$  and append token to the prefix
9: end for
10: Return:  $\mathbf{y}$  (accept if final state in  $\mathcal{F}$ )

```

---

## 6 Constrained decoding via grammars and finite-state constraints

Constrained decoding enforces a global constraint set  $\mathcal{C}(\mathbf{x})$  during generation. The key technical device is to represent  $\mathcal{C}(\mathbf{x})$  in a way that supports incremental feasibility checks on prefixes.

### 6.1 Prefix-closed feasibility and automata

A hard constraint set  $\mathcal{C} \subseteq \mathcal{V}^T$  is useful for incremental decoding when it induces a set of feasible prefixes

$$\text{Pref}(\mathcal{C}) := \{\mathbf{y}_{1:t} : \exists \mathbf{y}_{t+1:T} \text{ such that } (\mathbf{y}_{1:t}, \mathbf{y}_{t+1:T}) \in \mathcal{C}\}. \quad (13)$$

If a prefix  $\mathbf{y}_{1:t}$  is not in  $\text{Pref}(\mathcal{C})$ , then no continuation can satisfy the constraint, so it is safe to prune.

Finite-state constraints provide an especially clean representation. Let  $A = (\mathcal{S}, \mathcal{V}, \delta, s_0, \mathcal{F})$  be a finite-state automaton (FSA), where  $\mathcal{S}$  is a finite state set,  $s_0$  is the start state,  $\mathcal{F}$  is the accepting set, and  $\delta : \mathcal{S} \times \mathcal{V} \rightarrow 2^{\mathcal{S}}$  is the transition relation. A sequence  $\mathbf{y}$  is feasible if there exists a path from  $s_0$  to some  $s \in \mathcal{F}$  labeled by  $\mathbf{y}$ .

Given a current automaton state  $s_t$  after reading a prefix, the set of allowed next tokens is

$$\mathcal{A}(s_t) := \{v \in \mathcal{V} : \delta(s_t, v) \neq \emptyset\}. \quad (14)$$

A constrained next-token distribution is then obtained by masking and renormalizing:

$$\pi(v | y_{<t}, \mathbf{x}, s_t) := \frac{p_{\theta}(v | y_{<t}, \mathbf{x}) \mathbf{1}\{v \in \mathcal{A}(s_t)\}}{\sum_{v' \in \mathcal{A}(s_t)} p_{\theta}(v' | y_{<t}, \mathbf{x})}. \quad (15)$$

This is the simplest constrained decoding rule: it guarantees that generated prefixes remain feasible.

### 6.2 Grammar constraints

Grammars generalize FSAs. For regular constraints, an FSA is sufficient. For context-free constraints (for example matching brackets, structured queries, or typed expressions), one may use a parser state as the incremental constraint representation. The practical constraint-aware message is the same: at decoding step  $t$ , one computes a set of tokens that keep the partial parse valid, masks logits accordingly, and renormalizes.

In token-based LLMs, grammar constraints require care because tokenization breaks strings into subwords. A grammar written at the character level must be aligned with the token vocabulary, often by compiling the grammar into a token-level automaton. This compilation step is itself a constraint engineering decision, and it is one reason why constrained generation

often prefers domain-specific tokenizations (for example code tokens or JSON tokens) when strict syntactic validity is required.

### 6.3 Lexically constrained decoding

Lexical constraints require that a set of words or phrases appear in the generated sequence. This is common in translation (must include named entities), summarization (must include key facts), and domain generation (must include required fields). Lexical constraints can be represented by an automaton that tracks which constraints have been satisfied so far. If there are  $C$  required items, a conceptually simple representation is a state variable  $c \in \{0, 1\}^C$  that encodes which items have appeared.

Grid Beam Search (GBS) [Hokamp and Liu, 2017] constructs a grid over time  $t$  and the number of constraints satisfied, maintaining separate beams for each constraint-count cell and allowing only expansions that preserve feasibility. Dynamic Beam Allocation (DBA) [Post and Vilar, 2018] improves efficiency by allocating beam budget adaptively across constraint states rather than using a full grid.

**Example 1** (Lexical constraints as the first “policy” constraint). *Consider an email assistant that must generate a meeting confirmation containing two mandatory phrases: (i) “Rice Hall 340” and (ii) “Thursday at 9:30AM”. Let  $\mathcal{C}(\mathbf{x})$  be the set of continuations that contain both phrases at least once. A constrained decoder can maintain a state variable  $c \in \{0, 1\}^2$  and update it when a phrase is completed. Beam search can then be run over pairs  $(\text{prefix}, c)$ , pruning any expansion that makes it impossible to satisfy the remaining phrases within the length budget. This is a minimal instance of treating generation as search under a symbolic policy constraint.*

### 6.4 Complexity and failure modes

Constrained decoding trades probability mass for validity. If  $\mathcal{C}(\mathbf{x})$  is too strict relative to the base model, the feasible set may have extremely small probability, leading to brittle search and low-quality outputs. Algorithmically, constraints increase decoding complexity because the search state must track both the language model prefix and the constraint state (for example an automaton state or a constraint satisfaction vector). For beam-based methods, complexity typically scales as  $O(BT)$  in the unconstrained case, and as  $O(BT \cdot |\mathcal{S}|)$  or  $O(BT \cdot 2^C)$  depending on the constraint representation.

From the course perspective, this motivates three later themes. First, we want constraints to be represented in a way that admits efficient incremental feasibility checks. Second, we want the base model to assign nontrivial mass to the feasible set, which connects to training-time constraint injection. Third, we want principled ways to balance soft and hard constraints, which connects to optimization tools such as penalties, Lagrangians, and splitting methods.

## 7 Summary and outlook

Autoregressive Transformers define a tractable likelihood via the factorization  $p_{\theta}(\mathbf{y} \mid \mathbf{x}) = \prod_t p_{\theta}(y_t \mid y_{<t}, \mathbf{x})$  and implement it with a causal self-attention architecture. Decoding is then an inference procedure, either sampling or approximate optimization, and the decoding step is where constraints can be enforced directly as feasibility masks or as reranking energies.

The constrained decoding mechanisms in this lecture are intentionally simple: grammars, finite-state constraints, and lexical constraints. Their value is that they provide a clean, operational example of “generation as search subject to constraints” before we introduce the optimization toolbox. In later lectures, we will revisit the same pattern in continuous domains, where constraint enforcement often proceeds through projection, proximal operators, or guidance forces.

## References

- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020.
- A. Holtzman, J. Buys, L. Du, M. Forbes, and Y. Choi. The curious case of neural text degeneration. In *ICLR*, 2020.
- S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *NeurIPS*, 2015.
- M. Ranzato, S. Chopra, M. Auli, and W. Zaremba. Sequence level training with recurrent neural networks. In *ICLR*, 2016.
- L. Theis, A. van den Oord, and M. Bethge. A note on the evaluation of generative models. *arXiv:1511.01844*, 2016.
- C. Hokamp and Q. Liu. Lexically constrained decoding for sequence generation using grid beam search. In *ACL*, 2017.
- M. Post and D. Vilar. Fast lexically constrained decoding with dynamic beam allocation. In *NAACL*, 2018.
- Y. Li, Y. Liang, K. Zhang, and M. Ranzato. Contrastive decoding: Open-ended text generation as optimization. In *ACL*, 2022.
- A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are unsupervised multitask learners. Technical report, OpenAI, 2019.
- J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019.
- J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei. Scaling laws for neural language models. *arXiv:2001.08361*, 2020.
- J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas, L. Henderson, A. Menick, et al. Training compute-optimal large language models. *arXiv:2203.15556*, 2022.
- H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, et al. LLaMA: Open and efficient foundation language models. *arXiv:2302.13971*, 2023.
- J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, E. H. Chi, T. Hashimoto, P. Liang, and J. Dean. Emergent abilities of large language models. *arXiv:2206.07682*, 2022.