



X-Spanformer: A Tokenizer-Free, Span-Aware Encoder Inspired by X-Bar Theory

Kara Marie Rawson* Aimee Chrzanowski†

July 18, 2025

This work is a preprint v2 and has not yet been peer reviewed.

Abstract

Tokenization remains a fundamental bottleneck in transformer architectures, as static subword vocabularies fragment cross-domain text and impede real-time adaptability. We introduce X-Spanformer, a tokenizer-free, span-aware segmentation front end inspired by X-bar theory. X-Spanformer begins by inducing a hybrid Unigram-LM vocabulary—comprising all UTF-8 codepoints plus top entropy-pruned subword fragments—via an EM procedure approximated with Viterbi decoding, then applies a learned pruning criterion based on per-piece perplexity and OOV rate; this entire process is implemented as an efficient ONNX custom operator. The resulting soft piece-assignment matrix $P \in \mathbb{R}^{T \times V}$ is embedded and contextualized through a lightweight convolutional encoder that scales linearly in T . A factorized pointer network then identifies overlapping, variable-length spans, which are softly typed by modality, filtered by a learned length estimator, and pooled into d -dimensional span embeddings. We retain the top- K spans and fuse them via relevance-weighted interpolation into a single controller vector s , which is injected into downstream transformer layers through prefix-token, attention-bias, or gated fusion pathways. X-Spanformer is trained with an entropy-regularized span induction curriculum transitioning from synthetic supervision to type-aware signals. Empirical evaluation on compression ratio, contrastive retrieval alignment, span entropy, modality coherence, and end-to-end inference speed demonstrates that X-Spanformer outperforms static BPE and byte-level baselines in structural interpretability, compression efficiency, and throughput. An open-source ONNX-compatible implementation, along with training recipes and corpus construction guidelines, is provided to facilitate adoption across code, language, and hybrid domains.

1 Introduction

Transformer architectures underpin state-of-the-art performance in natural language understanding, code generation, and multimodal retrieval [4, 5, 6, 7]. A core prerequisite for these models is the segmentation of raw text or code into discrete units, most often fixed subwords derived from Byte-Pair Encoding (BPE) [8] or Unigram-LM vocabularies [1]. Although static vocabularies enable efficient embedding lookup and in-domain accuracy, they introduce rigid lexical boundaries that (i) degrade performance under domain shift [9], (ii) obscure long-range compositional patterns in code,

*rawsonkara@gmail.com

†aimeechrzanowski@gmail.com

multilingual, and hybrid text, and (iii) require costly re-training or vocabulary expansion to handle novel syntax or semantics. Moreover, treating segmentation as an irreversible preprocessing step prevents gradient flow and adaptation to downstream objectives.

To address these limitations, we propose X-Spanformer, a tokenizer-free, span-aware segmentation module that is fully differentiable and ONNX-native. X-Spanformer replaces static tokenization with an adaptive, corpus-driven induction of a hybrid Unigram-LM vocabulary via an expectation–maximization routine approximated by Viterbi decoding, followed by entropy- and perplexity-driven pruning to control vocabulary size and coverage. The resulting soft assignment matrix $P \in \mathbb{R}^{T \times V}$ is projected into seed embeddings and contextualized by a lightweight convolutional encoder that scales linearly in sequence length. A factorized pointer network locates overlapping, variable-length span candidates, which are softly typed by modality entropy, filtered by a learned length estimator, and pooled into d -dimensional span embeddings. We fuse the top- K spans into a single controller vector via relevance-weighted interpolation and inject it into downstream transformers through bias, prefix, or gated-FFN mechanisms. Because every component is differentiable and ONNX-compatible, X-Spanformer supports end-to-end training with an entropy-regularized span induction curriculum, yielding interpretable segmentation aligned with phrase-level structure and efficient inference across domains.

1.1 Contributions

This paper makes the following contributions:

1. We formalize tokenizer-free segmentation as a span-prediction problem grounded in X-bar theory, driven by a hybrid Unigram-LM front end with Viterbi-based EM and adaptive perplexity/OOV pruning.
2. We introduce a convolutional contextual encoder that augments seed embeddings with local structure at linear cost, prioritizing throughput and interpretability over global self-attention.
3. We design an ONNX-native architecture for span pooling, modality-aware scoring, and controller fusion, detailing multiple injection pathways (prefix, bias, gated-FFN) for flexible downstream integration.
4. We develop an entropy-regularized, multi-phase curriculum that bootstraps from synthetic BPE supervision to type-aware signals, enabling robust span induction across code, language, and hybrid modalities.
5. We propose a comprehensive evaluation suite—including compression ratio, contrastive retrieval alignment, span entropy, modality coherence, and inference latency—and release an open-source ONNX implementation with training recipes and corpus construction guidelines.

2 Related Work

2.1 Static Subword Tokenization

Conventional transformer pipelines rely on fixed subword vocabularies, most prominently Byte-Pair Encoding (BPE) [8] and Unigram-LM tokenizers such as SentencePiece [1]. These methods facilitate

efficient embedding lookup and strong in-domain performance, but their immutable vocabularies cannot adapt during training, fragment compositional or rare phrases under domain shift [9], and necessitate costly re-tokenization when new syntax or terminology arise.

2.2 Character- and Byte-Level Encoders

To eliminate subword heuristics, character- and byte-level models process raw codepoints directly. Charformer [10] learns latent splits via gradient-based tokenization, CANINE [11] applies down- and up-sampling around raw Unicode, and ByT5 [15] explores purely byte-level inputs. While these approaches remove offline segmentation, they do not produce overlapping or hierarchical spans and rely on global self-attention for context.

2.3 Differentiable Segmentation

Unsupervised and differentiable segmentation integrates boundary induction into neural pretraining. Morfessor [16] applies minimum description-length to discover morphemes. Probabilistically Masked Language Models [13] and learned-segmentation models [12] optimize reconstruction objectives to induce non-overlapping partitions. These methods enable end-to-end learning but lack explicit modeling of overlapping spans or linguistic structure.

2.4 Span-Based and Pointer Networks

Pointer networks [3] predict arbitrary start–end pairs to locate variable-length spans. SpanBERT [17] leverages this for masked span reconstruction, and similar architectures in vision and speech generate overlapping proposals for improved alignment [18, 19]. However, prior work has not combined pointer-based, overlapping span induction with explicit grammar-inspired priors within transformer encoders.

2.5 Bridging Static and Dynamic Segmentation

Hybrid vocabularies that merge character coverage with static subwords have been proposed to improve robustness [10, 11]. Our approach extends this paradigm by training a hybrid Unigram-LM front end—induced by an EM procedure approximated via Viterbi decoding and pruned on perplexity and OOV thresholds—implemented as an ONNX custom operator. We then employ a convolutional encoder to contextualize soft assignments, a factorized pointer network to propose overlapping spans, and a modality-aware scoring mechanism to filter and fuse spans into a controller vector for downstream transformer injection.

In contrast to purely static tokenizers, character-level models, and non-overlapping segmentation methods, X-Spanformer unifies hybrid vocabulary induction, differentiable Viterbi-based pruning, convolutional contextualization, and pointer-network span prediction under a single, ONNX-native framework. This design yields dynamic, interpretable spans that adapt to domain shifts and can be optimized end-to-end alongside downstream tasks.

3 Architecture

This section formalizes the modular components of X-Spanformer and their interactions within the segmentation pipeline. Each architectural unit is motivated, given a precise mathematical formulation, and illustrated with pseudocode where appropriate. We conclude with strategies for integrating the fused span controller into standard transformer encoders and analyze runtime complexity.

X-Spanformer begins with a hybrid Unigram-LM vocabulary¹ that replaces hard token boundaries with soft, probabilistic segmentation. A custom ONNX operator ingests raw Unicode codepoints of length T and computes a probability matrix

$$P \in \mathbb{R}^{T \times V}, \quad P_{t,i} = \Pr(\text{piece } u_i \text{ starts at position } t),$$

where V is the induced vocabulary size. We embed these soft piece probabilities via

$$H^0 = P W_{\text{emb}}, \quad W_{\text{emb}} \in \mathbb{R}^{V \times d},$$

yielding seed embeddings $H^0 \in \mathbb{R}^{T \times d}$. From these low-level vectors, the model:

- Extracts a ranked span set $S = \{(i_k, j_k)\}_{k=1}^K$ with $1 \leq i_k < j_k \leq T$ via a pointer network [3].
- Pools each span (i_k, j_k) into a d -dimensional embedding $s_{i_k j_k}$ (mean or gated pooling) [10].
- Predicts soft modality distributions $p_{i_k j_k}^{\text{type}} \in \Delta^M$, reflecting types such as code, natural language, or identifier [20, 21].
- Filters down to a final span set $S' \subseteq S$ via a learned length estimator [22].

Span-level augmentation parallels auxiliary token insertion in models like SpanBERT [17] but operates on soft, overlapping spans rather than hard subwords. All modules—from boundary scoring to controller fusion—are fully differentiable and trained end-to-end with the downstream encoder.

3.1 Hybrid Unigram-LM Vocabulary Induction

We construct a data-adaptive subword vocabulary by training a Unigram language model (LM) over all substrings up to length L_{\max} , then pruning pieces based on both corpus perplexity and coverage (OOV rate). Let $\mathcal{X} = \{x^{(i)}\}_{i=1}^N$ be our corpus of raw Unicode codepoint sequences of total length T , and let \mathcal{U}_0 denote the initial candidate set formed by extracting every substring of length $\leq L_{\max}$, retaining the top M by frequency, and including all individual codepoints.

3.1.1 Candidate Generation

For each sequence $x \in \mathcal{X}$, we slide a window of width $\leq L_{\max}$ across all codepoints to extract substrings. We record their frequency and retain the top M most frequent substrings. The resulting candidate set is:

$$\mathcal{U}_0 = \{\text{top } M \text{ substrings}\} \cup \{\text{all single codepoints}\},$$

ensuring base coverage and bounding vocabulary size.

¹Learned via EM over substrings up to length L_{\max} ; see Sec. 3.1 for details.

3.1.2 Unigram LM via EM

We initialize the piece probabilities by normalizing raw frequency counts:

$$p^{(0)}(u) = \frac{\text{freq}(u)}{\sum_{v \in \mathcal{U}_0} \text{freq}(v)}, \quad \forall u \in \mathcal{U}_0.$$

At each EM iteration t , we alternate:

E-step: For each sequence $x \in \mathcal{X}$, we compute the posterior usage mass $\gamma^{(t)}(u | x)$ for every piece $u \in \mathcal{U}_t$. Since enumerating all segmentations of x is intractable, we approximate using the single best segmentation:

$$\text{seg}^*(x) = \arg \max_{\text{seg}} \prod_{v \in \text{seg}} p^{(t)}(v),$$

obtained via Viterbi decoding [1]. We then compute:

$$\gamma^{(t)}(u | x) = \sum_{v \in \text{seg}^*(x)} \mathbf{1}_{u=v}.$$

M-step: We re-estimate the piece probabilities as:

$$p^{(t+1)}(u) = \frac{\sum_{x \in \mathcal{X}} \gamma^{(t)}(u | x)}{\sum_{v \in \mathcal{U}_t} \sum_{x \in \mathcal{X}} \gamma^{(t)}(v | x)},$$

which reduces to relative frequency over the Viterbi segmentations.

3.1.3 Baseline Perplexity

To benchmark compressibility before pruning, we compute the initial corpus perplexity:

$$\text{PPL}^{(0)} = \exp \left(-\frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} \log \prod_{v \in \text{seg}^*(x)} p^{(0)}(v) \right),$$

where $\text{seg}^*(x)$ is the Viterbi segmentation under $p^{(0)}$.

3.1.4 Adaptive Pruning by PPL and OOV

After each M-step, we consider pruning any piece $u \in V$ with $p^{(t)}(u) < \epsilon$. For a candidate removal $V' = V \setminus \{u\}$, we first perform Viterbi segmentation under the reduced vocabulary:

$$\text{seg}_{V'}^*(x) = \arg \max_{\text{seg}} \prod_{v \in \text{seg}} p_{V'}^{(t)}(v) \quad \text{for each } x \in \mathcal{X}.$$

We then introduce the following quantities:

$$N'_p = \sum_{x \in \mathcal{X}} |\text{seg}_{V'}^*(x)| \quad (\text{total number of subword pieces}),$$

$$L' = - \sum_{x \in \mathcal{X}} \sum_{v \in \text{seg}_{V'}^*(x)} \log p_{V'}^{(t)}(v) \quad (\text{negative log-likelihood of all pieces}),$$

$$N_t = \sum_{x \in \mathcal{X}} |x| \quad (\text{total number of codepoint tokens}),$$

$$N'_{\text{uncov}} = \sum_{x \in \mathcal{X}} \sum_{i=1}^{|x|} \mathbf{1}(i \notin \text{cover}_{V'}(x)) \quad (\text{total uncovered positions}),$$

where $\text{cover}_{V'}(x)$ is the set of codepoint indices spanned by $\text{seg}_{V'}^*(x)$.

Using these, we define:

$$\text{PPL}' = \exp\left(\frac{L'}{N'_p}\right), \quad \text{OOV}' = \frac{N'_{\text{uncov}}}{N_t}.$$

We accept the removal u if and only if both

$$\text{PPL}' - \text{PPL}^{(t)} < \tau_{\text{ppl}} \quad \text{and} \quad \text{OOV}' \leq \delta_{\text{oov}}.$$

This criterion guarantees that any pruning step degrades the model's average piece-level log-likelihood by at most τ_{ppl} and introduces no more than δ_{oov} uncovered codepoint positions.

3.1.5 Existence and Monotonicity Proposition

Define the feasible set:

$$\mathcal{F}(\tau, \delta) = \left\{ V \subseteq \mathcal{U}_0 \mid \text{PPL}(V) \leq \text{PPL}^{(0)} + \tau, \text{ OOV}(V) \leq \delta \right\}.$$

Proposition 1 (Feasibility and Monotonicity). *For any $\tau, \delta \geq 0$, the feasible set $\mathcal{F}(\tau, \delta)$ is nonempty. Moreover, if $\tau' \geq \tau$, $\delta' \geq \delta$, then:*

$$\mathcal{F}(\tau, \delta) \subseteq \mathcal{F}(\tau', \delta'), \quad \Rightarrow \quad \min_{V \in \mathcal{F}(\tau', \delta')} |V| \leq \min_{V \in \mathcal{F}(\tau, \delta)} |V|.$$

Proof. **Step 1: Existence**

Choose $V = \mathcal{U}_0$, the full candidate set.

Since V includes all substrings and individual codepoints:

$$\text{OOV}(V) = 0, \quad \text{PPL}(V) = \text{PPL}^{(0)}.$$

Therefore, for any $\tau, \delta \geq 0$,

$$\text{PPL}(V) \leq \text{PPL}^{(0)} + \tau, \quad \text{OOV}(V) \leq \delta, \quad \Rightarrow V \in \mathcal{F}(\tau, \delta).$$

Step 2: Monotonicity

Let $V \in \mathcal{F}(\tau, \delta)$, and suppose $\tau' \geq \tau, \delta' \geq \delta$. Then:

$$\text{PPL}(V) \leq \text{PPL}^{(0)} + \tau \leq \text{PPL}^{(0)} + \tau', \quad \text{OOV}(V) \leq \delta \leq \delta', \quad \Rightarrow V \in \mathcal{F}(\tau', \delta').$$

Thus:

$$\mathcal{F}(\tau, \delta) \subseteq \mathcal{F}(\tau', \delta').$$

Step 3: Minimality

Since $\mathcal{F}(\tau, \delta) \subseteq \mathcal{F}(\tau', \delta')$,

$$\min_{V \in \mathcal{F}(\tau', \delta')} |V| \leq \min_{V \in \mathcal{F}(\tau, \delta)} |V|.$$

□

3.1.6 Algorithm: Adaptive Vocabulary Induction

To construct an efficient and corpus-aware subword vocabulary, we apply an expectation-maximization loop over the candidate set \mathcal{U}_0 , guided by a Viterbi-decoded likelihood objective. Starting from raw frequency estimates, we refine piece probabilities $p(u)$ via best-segmentation token counts and normalize them across the corpus. After each M-step, we attempt to prune low-probability pieces while preserving compressibility and coverage. Each candidate removal is simulated by resegmenting \mathcal{X} under the reduced vocabulary V' and evaluating both perplexity and position-level OOV rate. If the degradation in log-likelihood is bounded and no coverage gaps are introduced, the removal is accepted. This selective pruning yields a compact, entropy-regularized vocabulary tailored to the domain structure.

Algorithm 1 Adaptive Unigram-LM Vocabulary Induction

```

1: Extract candidate substrings up to length  $L_{\max}$  from corpus  $\mathcal{X}$ ; form initial vocabulary  $\mathcal{U}_0$ 
2: Initialize piece probabilities:  $p^{(0)}(u) \propto \text{freq}(u)$  for all  $u \in \mathcal{U}_0$ 
3: Compute baseline perplexity  $\text{PPL}^{(0)}$  via Viterbi decoding over  $\mathcal{X}$ 
4: for iteration  $t = 0$  to  $T_{\max}$  do
5:   for each sequence  $x \in \mathcal{X}$  do
6:     Compute best segmentation  $\text{seg}^*(x)$  using current  $p^{(t)}$ 
7:     Accumulate token usage counts  $\gamma^{(t)}(u | x)$  for all  $u \in \text{seg}^*(x)$ 
8:   end for
9:   Update probabilities: normalize counts to get  $p^{(t+1)}(u)$ 
10:  Set current vocabulary  $V = \{u \mid p^{(t+1)}(u) > 0\}$ 
11:  for each  $u \in V$  with  $p^{(t+1)}(u) < \epsilon$  do
12:    Tentatively prune: define  $V' = V \setminus \{u\}$ 
13:    for each sequence  $x \in \mathcal{X}$  do
14:      Decode  $\text{seg}_{V'}^*(x)$  using updated  $V'$ 
15:      Compute log-prob score and uncovered positions
16:    end for
17:    Compute  $\text{PPL}'$  and  $\text{OOV}'$  from new segmentations
18:    if  $\text{PPL}' - \text{PPL}^{(t)} < \tau_{\text{ppl}}$  and  $\text{OOV}' \leq \delta_{\text{oov}}$  then
19:      Accept removal:  $V \leftarrow V'$ 
20:    end if
21:  end for
22:  Update  $\text{PPL}^{(t+1)}$  using accepted vocabulary  $V$ 
23: end for
24: Return final pruned vocabulary  $V$  and piece probabilities  $\{p(u)\}$ 

```

3.2 Seed Embeddings and Candidate Set

The X-Spanformer pipeline begins with a sequence of T raw Unicode codepoints. We apply our custom Unigram-LM operator (Sec. 3.1), implemented as an ONNX custom op, to compute soft piece probabilities over the induced vocabulary V :

$$P \in \mathbb{R}^{T \times V}, \quad P_{t,i} = \Pr(\text{piece } u_i \text{ starts at position } t).$$

The vocabulary V comprises all single codepoints plus the final entropy- and perplexity-pruned subword fragments.

These probabilities are embedded via a single matrix multiplication:

$$H^0 = P W_{\text{emb}}, \quad W_{\text{emb}} \in \mathbb{R}^{V \times d},$$

yielding initial seed embeddings $H^0 \in \mathbb{R}^{T \times d}$ for each input position.

3.2.1 Contextual Encoder: Convolutional Pathway

To enrich these seed embeddings with contextual structure, we apply a lightweight convolutional encoder to produce the final representation $H \in \mathbb{R}^{T \times d}$:

$$H = \text{ConvNet}(H^0), \quad \text{cost } O(Td^2).$$

This encoder employs stacked or dilated 1D convolutions [10], providing fixed-window receptive fields that are tuned to capture short-range compositional patterns. By avoiding the quadratic complexity of self-attention, the convolutional encoder ensures scalability across long codepoint sequences and maintains a runtime that is linear in the input length.

The design choice to adopt convolutional contextualization reflects the architectural goal of high-throughput segmentation in real-time or resource-constrained settings. While transformer-based encoders offer richer modeling of long-distance dependencies [4], their runtime limits applicability for streaming inputs or modular deployments. In contrast, convolutional filters preserve local positional regularities critical for span induction—such as those arising in identifiers, literals, or markup—and enable reproducible, interpretable segmentation across diverse modalities. Empirical evidence on synthetic benchmarks further suggests that this encoder achieves robust span localization without sacrificing structural fidelity [2].

3.2.2 Span Candidate Enumeration

Given the contextualized sequence H , we enumerate all contiguous span candidates:

$$C = \{(i, j) \mid 1 \leq i < j \leq T, j - i + 1 \leq w_{\max}\}.$$

This restricts span width to w_{\max} , reducing total candidates from quadratic to linear scale: $|C| = O(T w_{\max})$. Each span (i, j) corresponds to the subsequence $[h_i, \dots, h_j]$, which is passed to the boundary scorer and length filter in subsequent modules. The candidate set remains compatible with global-attention filtering and sparse attention schemes [17, 10], ensuring scalability and modularity in downstream processing.

3.3 Span Predictor

Given contextualized embeddings

$$H \in \mathbb{R}^{T \times d},$$

we first form the candidate span set:

$$C = \{(i, j) \mid 1 \leq i < j \leq T, j - i \leq w_{\max}\},$$

as defined in Sec. 3.2. We then use two parallel linear heads—forming a factorized pointer network [3]—to predict span boundaries:

$$\ell^s = W_s H + b_s, \quad p^s = \text{softmax}(\ell^s), \quad \ell^e = W_e H + b_e, \quad p^e = \text{softmax}(\ell^e),$$

where $W_s, W_e \in \mathbb{R}^{T \times d}$, $b_s, b_e \in \mathbb{R}^T$, and p_i^s, p_j^e denote the probabilities that a span begins at position i and ends at position j , respectively.

Each span $(i, j) \in C$ is scored by the outer-product of its boundary probabilities:

$$\text{score}(i, j) = p_i^s p_j^e.$$

This method efficiently captures boundary salience and mirrors successful strategies in QA and span-based extraction models [23, 24].

We then select the top- K spans by score:

$$S = \text{TopK}\{\text{score}(i, j) \mid (i, j) \in C\}.$$

Proposition 2 (Top- K Marginal Likelihood). *Let $p^s, p^e \in \Delta^T$ be independent distributions over start and end positions. For each span $(i, j) \in C$, define*

$$P(i, j) = p_i^s p_j^e.$$

Then the top- K spans by $P(i, j)$ maximize total mass among all subsets of size K :

$$S = \arg \max_{\substack{S' \subseteq C \\ |S'|=K}} \sum_{(i,j) \in S'} P(i, j).$$

Proof. We aim to solve

$$\max_{\substack{S' \subseteq C \\ |S'|=K}} \sum_{(i,j) \in S'} P(i, j).$$

Step 1: The objective is additive in $P(i, j)$, and all terms are non-negative:

$$P(i, j) \geq 0, \quad \text{so} \quad \sum_{(i,j) \in S'} P(i, j) \text{ increases with high-mass elements.}$$

Step 2: Sorting all $(i, j) \in C$ by $P(i, j)$, selecting the K largest entries yields the subset with maximal total mass:

$$S = \text{TopK}\{P(i, j)\}.$$

Step 3: Since p^s and p^e are independent, there are no interaction terms across spans—greedy selection remains optimal.

Conclusion:

$$S = \arg \max_{|S'|=K} \sum_{(i,j) \in S'} p_i^s p_j^e$$

as claimed. \square

3.4 Length Estimator

Even high-confidence boundary proposals can yield spans that are too short or too long to reflect meaningful linguistic units. To inject a learned prior over span width, we train a categorical length estimator that filters candidates based on predicted versus actual length.

Span Length and Pooling For each candidate span $(i, j) \in S$, define the true length:

$$\delta = j - i + 1.$$

We pool the contextual embeddings over the span window:

$$v_{ij} = \text{Pool}(H[i:j]) \in \mathbb{R}^d,$$

where $\text{Pool}(\cdot)$ may be mean-, max-, gated-, or self-attentive aggregation [10].

Length Prediction We predict a categorical distribution over B discrete length bins:

$$\ell^\delta = W_\ell v_{ij} + b_\ell, \quad p^\delta = \text{softmax}(\ell^\delta), \quad \hat{\delta} = \arg \max p^\delta,$$

where $\hat{\delta}$ serves as a learned prior on plausible span width.

Tolerance-Based Filtering We retain only spans whose true length δ lies within a tolerance τ of the predicted bin:

$$S' = \{(i, j) \in S \mid |(j - i + 1) - \hat{\delta}| \leq \tau\}.$$

The hyperparameter τ determines how strictly the length prediction must match the actual span width.

Proposition 3 (Span Count Upper Bound). *Assume all learned lengths $\hat{\delta} \in [\delta_{\min}, \delta_{\max}]$ and fix a tolerance $\tau < \delta_{\max} - \delta_{\min}$. Then the total number of retained spans satisfies:*

$$|S'| = \mathcal{O}(T \cdot (2\tau + 1)),$$

i.e., only $O(T)$ spans are retained per input sequence.

Proof. We analyze the filtering region induced by length tolerance:

1. For a fixed start position i , the predicted length is $\hat{\delta}$.

2. The allowable end index j satisfies:

$$j = i + \hat{\delta} - 1 \pm \tau \Rightarrow j \in [i + \hat{\delta} - \tau - 1, i + \hat{\delta} + \tau - 1].$$

3. Thus, for each i , the number of valid j positions is at most:

$$(i + \hat{\delta} + \tau - 1) - (i + \hat{\delta} - \tau - 1) + 1 = 2\tau + 1.$$

4. There are T possible start positions i , so:

$$|S'| \leq T \cdot (2\tau + 1) = \mathcal{O}(T).$$

□

This learned filtering module reduces the subquadratic span candidate space to linear size, while enforcing cognitively motivated length regularity [14]. The tolerance parameter τ controls the trade-off between structural fidelity and coverage, improving both efficiency and interpretability.

3.5 Modality Typing

Text streams often mix multiple domains—natural language, code syntax, identifiers, numeric values, markup—especially in technical or hybrid contexts. To model this structure, we attach a modality classifier to each pooled span embedding $v_{ij} \in \mathbb{R}^d$, as defined in Sec. 3.4. The classifier predicts a soft type distribution and accompanying entropy score.

3.5.1 Modality Classification

Let M be the number of modality classes. We first transform the span vector:

$$h_{ij} = \text{ReLU}(W_v v_{ij} + b_v) \in \mathbb{R}^{d'},$$

where $d' \leq d$. We then compute:

$$\ell_{ij}^{\text{mod}} = W_{\text{mod}} h_{ij} + b_{\text{mod}}, \quad p_{ij}^{\text{mod}} = \text{softmax}(\ell_{ij}^{\text{mod}}) \in \Delta^M.$$

Here $p_{ij,m}^{\text{mod}}$ is the model’s belief that span (i, j) belongs to modality m .

3.5.2 Modality Entropy and Ambiguity

To quantify ambiguity:

$$H_{ij}^{\text{mod}} = - \sum_{m=1}^M p_{ij,m}^{\text{mod}} \log p_{ij,m}^{\text{mod}}.$$

Higher entropy indicates uncertainty, useful for early training exploration or hybrid domains.

3.5.3 Auxiliary Supervision

When gold labels $y_{ij}^{\text{gold}} \in \{0, 1\}^M$ are available, we optimize:

$$\mathcal{L}_{\text{mod}} = - \sum_{(i,j) \in S} \sum_{m=1}^M y_{ij,m}^{\text{gold}} \log p_{ij,m}^{\text{mod}}.$$

This improves zero-shot transfer across modalities [25, 26].

3.5.4 Conditional Routing

During inference, p_{ij}^{mod} can modulate decoder behavior:

$$e_{ij}^{\text{mod}} = p_{ij}^{\text{mod}} \cdot E_{\text{mod}}, \quad E_{\text{mod}} \in \mathbb{R}^{M \times d}.$$

This embedding may be concatenated to v_{ij} or used to bias attention heads [21].

3.5.5 Interpretability

The modality entropy and distribution expose latent semantic domains of each span, aiding mixed-modality diagnostics and span alignment analysis [20, 10].

3.5.6 Integration into Span Scoring

We incorporate modality typing into the span relevance MLP f_{score} (Sec. 3.8). Define:

$$d_{ij} = j - i + 1, \quad c_{ij} = p_i^s p_j^e.$$

Let $\phi(d_{ij}) \in \mathbb{R}^D$ be a learned embedding of span length. Form the joint feature vector:

$$x_{ij} = [v_{ij}; \phi(d_{ij}); p_{ij}^{\text{mod}}; H_{ij}^{\text{mod}}; c_{ij}] \in \mathbb{R}^{d+D+M+1+1}.$$

We compute relevance weight:

$$w_{ij} = \text{MLP}_{\text{score}}(x_{ij}), \quad a_{ij} = \frac{\exp(w_{ij})}{\sum_{(p,q) \in S'} \exp(w_{pq})}.$$

Proposition 4 (Modality-Aware Relevance Weighting). *Let $x_{ij} \in \mathbb{R}^d$ be a span descriptor including content, length, boundary confidence, and modality entropy. Then the normalized weight*

$$a_{ij} = \frac{\exp(\text{MLP}(x_{ij}))}{\sum_{(p,q) \in S'} \exp(\text{MLP}(x_{pq}))}$$

defines a probability distribution over spans in S' , with preference toward spans that match learned structural and semantic type priors.

Proof. Step 1: Positivity and normalization

For each $(i,j) \in S'$, let $w_{ij} = \text{MLP}(x_{ij})$. Since $\exp(w_{ij}) > 0$, it follows immediately that

$$a_{ij} > 0.$$

Moreover,

$$\sum_{(i,j) \in S'} a_{ij} = \sum_{(i,j) \in S'} \frac{\exp(w_{ij})}{\sum_{(p,q) \in S'} \exp(w_{pq})} = \frac{\sum_{(i,j) \in S'} \exp(w_{ij})}{\sum_{(p,q) \in S'} \exp(w_{pq})} = 1.$$

Step 2: Preference via logits

Because a_{ij} is obtained by applying the softmax to the logits $\{w_{pq}\}$, it is strictly increasing in its own logit and decreasing in the others. Concretely, if for two spans $w_{ij} > w_{p'q'}$, then

$$a_{ij} = \frac{\exp(w_{ij})}{\sum \exp(w)} > \frac{\exp(w_{p'q'})}{\sum \exp(w)} = a_{p'q'}.$$

Since $w_{ij} = \text{MLP}(x_{ij})$ incorporates modality distribution p_{ij}^{mod} and entropy H_{ij}^{mod} , spans with lower entropy and stronger modality alignment receive higher w_{ij} and thus larger a_{ij} .

Step 3: Explicit softmax form

Combining the above, we recover the stated formula

$$a_{ij} = \frac{\exp(\text{MLP}(x_{ij}))}{\sum_{(p,q) \in S'} \exp(\text{MLP}(x_{pq}))},$$

which is a valid probability distribution over S' that embeds learned modality-aware preferences. \square

3.6 Span Embedding

Each retained span $(i,j) \in S'$ is mapped to a fixed-length vector $s_{ij} \in \mathbb{R}^d$ that encodes both its internal composition and contextual salience. To achieve this, we employ a dual encoder: a lightweight mean-pooling pathway and a local self-attention pathway, fused adaptively via a gated interpolator.

3.6.1 Mean-Pooling Encoder

Define the span length $\delta = j - i + 1$ and compute the average of its token embeddings:

$$\bar{h}_{ij} = \frac{1}{\delta} \sum_{k=i}^j h_k.$$

Project the pooled vector into model dimension:

$$s_{ij}^{\text{mean}} = W_m \bar{h}_{ij} + b_m, \quad W_m \in \mathbb{R}^{d \times d}.$$

Mean pooling is efficient and length-invariant, and performs well in extractive tasks [27, 17].

3.6.2 Local Self-Attention Encoder

We apply multi-head self-attention over the span window $H[i:j]$:

$$\{A_{ij}^{(\ell)}\}_{\ell=1}^h = \text{MHSA}(H[i:j]), \quad A_{ij} = [A_{ij}^{(1)}; \dots; A_{ij}^{(h)}].$$

Then project the concatenated heads:

$$s_{ij}^{\text{attn}} = W_o A_{ij} + b_o, \quad W_o \in \mathbb{R}^{d \times (h d_h)},$$

where $d_h = d/h$ and h is the number of heads. This encoder captures intra-span dependencies and asymmetries [28, 10].

3.6.3 Gated Fusion

To balance expressivity with compute, we learn a gate $g_{ij} \in (0, 1)$ that weights the two embeddings. Define the fusion feature vector:

$$f_{ij} = [\bar{h}_{ij}; \phi(\delta); p_{ij}^{\text{mod}}; H_{ij}^{\text{mod}}; c_{ij}] \in \mathbb{R}^{d+D+M+1+1},$$

where: - $\phi(\delta) \in \mathbb{R}^D$ is a learned length embedding, - $p_{ij}^{\text{mod}} \in \Delta^M$ is the modality distribution (Sec. 3.5), - $H_{ij}^{\text{mod}} \in [0, \log M]$ is modality entropy, - $c_{ij} = p_i^s p_j^e$ is boundary confidence.

Then:

$$g_{ij} = \sigma(w_g^\top f_{ij} + b_g), \quad s_{ij} = g_{ij} s_{ij}^{\text{attn}} + (1 - g_{ij}) s_{ij}^{\text{mean}}.$$

Proposition 5 (Adaptive Fusion Bound). *Let both $s_{ij}^{\text{mean}}, s_{ij}^{\text{attn}} \in \mathbb{R}^d$ have bounded norm $\|s\| \leq B$, and $g_{ij} \in [0, 1]$. Then:*

$$\|s_{ij}\|_2 \leq B.$$

Proof. We compute:

$$\|s_{ij}\|_2 = \|g_{ij} s_{ij}^{\text{attn}} + (1 - g_{ij}) s_{ij}^{\text{mean}}\|_2 \leq g_{ij} \|s_{ij}^{\text{attn}}\|_2 + (1 - g_{ij}) \|s_{ij}^{\text{mean}}\|_2 \leq g_{ij} B + (1 - g_{ij}) B = B.$$

This uses convexity and the triangle inequality. \square

All final span embeddings $s_{ij} \in \mathbb{R}^d$ are forwarded to the controller-fusion stage (Sec. 3.7) and participate in downstream interpolation, type prediction, and transformer injection.

3.7 Controller Fusion

Given the filtered set of span embeddings $\{s_{ij} \in \mathbb{R}^d\}_{(i,j) \in S'}$, we compute a single fused controller vector $s \in \mathbb{R}^d$ via relevance-weighted interpolation. Each span is scored based on its content, type, structure, and boundary confidence.

3.7.1 Feature Construction and Scoring

Define the relevance feature for span $(i, j) \in S'$:

$$x_{ij} = [s_{ij}; \phi(\delta_{ij}); p_{ij}^{\text{mod}}; H_{ij}^{\text{mod}}; c_{ij}] \in \mathbb{R}^{d+D+M+1+1},$$

where: - $\delta_{ij} = j - i + 1$ is the span length, - $\phi(\delta_{ij}) \in \mathbb{R}^D$ is a learned embedding, - $p_{ij}^{\text{mod}} \in \Delta^M$ is the modality distribution, - H_{ij}^{mod} is the entropy, - $c_{ij} = p_i^s p_j^e$ is boundary confidence.

Compute normalized relevance weights:

$$w_{ij} = f_{\text{score}}(x_{ij}), \quad a_{ij} = \frac{\exp(w_{ij})}{\sum_{(p,q) \in S'} \exp(w_{pq})}.$$

Fuse the controller vector:

$$s = \sum_{(i,j) \in S'} a_{ij} \cdot s_{ij}.$$

Proposition 6 (Controller Interpolation Properties). *Let each $s_{ij} \in \mathbb{R}^d$ have bounded norm $\|s_{ij}\| \leq B$, and weights $a_{ij} \in [0, 1]$, $\sum a_{ij} = 1$. Then:*

$$s \in \text{conv}(\{s_{ij}\}), \quad \|s\|_2 \leq B.$$

Proof. Since a_{ij} form a probability distribution and all inputs lie in a bounded ball of radius B , their convex combination also lies within that ball:

$$\|s\|_2 = \left\| \sum a_{ij} s_{ij} \right\|_2 \leq \sum a_{ij} \|s_{ij}\|_2 \leq \sum a_{ij} B = B.$$

□

3.7.2 Transformer Injection Modes

The controller vector s is injected into the transformer encoder $\text{Transf}(\cdot)$ via three differentiable modes:

Prefix-Token Injection Prepend s as a synthetic token:

$$H' = [s; h_1; \dots; h_T] \in \mathbb{R}^{(T+1) \times d},$$

with position embedding for the prefix. This enables global access from layer 0 [21].

Attention-Bias Injection Bias queries and keys:

$$Q_i \leftarrow Q_i + W_Q s, \quad K_j \leftarrow K_j + W_K s,$$

so attention scores become:

$$A_{ij} = \exp((Q_i + W_Q s)^\top (K_j + W_K s)).$$

Here $W_Q, W_K \in \mathbb{R}^{d \times d}$ control how s shifts attention geometry [30].

Gated-FFN Injection Modulate feed-forward response:

$$g = \sigma(W_g s + b_g) \in (0, 1)^d, \quad h'_i = h_i + g \odot \text{FFN}(h_i),$$

where \odot denotes elementwise multiplication. This adaptively gates nonlinear transformation based on span context [31].

3.7.3 Multi-Path Fusion

We learn nonnegative scalar weights α, β, γ to interpolate injection modes:

$$\text{Transf}(\alpha(\text{PREFIX}) + \beta(\text{ATTN-BIAS}) + \gamma(\text{GATED-FFN})),$$

preserving full differentiability and allowing dynamic tuning of injection strategy.

3.8 Span Interpolation for Overlap Resolution

After filtering, we obtain a set of overlapping span embeddings $\{s_{ij}\}_{(i,j) \in S'} \subset \mathbb{R}^d$. Rather than injecting each span separately, we compute a single fused controller vector $\tilde{s} \in \mathbb{R}^d$ via relevance-weighted interpolation.

3.8.1 Soft Interpolation

For each span $(i, j) \in S'$, we define its relevance logit:

$$w_{ij} = f_{\text{score}}(s_{ij}, \phi(\delta_{ij}), p_{ij}^{\text{mod}}, H_{ij}^{\text{mod}}, \log c_{ij}), \quad (1)$$

where: - $\delta_{ij} = j - i + 1$, - $\phi(\delta_{ij}) \in \mathbb{R}^D$ is a learned embedding, - $p_{ij}^{\text{mod}} \in \Delta^M$ is the modality distribution, - H_{ij}^{mod} is entropy (Sec. 3.5), - $c_{ij} = p_i^s p_j^e$ is the boundary confidence.

We normalize via softmax:

$$\alpha_{ij} = \frac{\exp(w_{ij})}{\sum_{(p,q) \in S'} \exp(w_{pq})}, \quad \tilde{s} = \sum_{(i,j) \in S'} \alpha_{ij} s_{ij}. \quad (2)$$

This formulation parallels soft-attention in retrieval-augmented models [32, 33] and expert fusion [34].

3.8.2 Theoretical Properties

Proposition 7 (Span Interpolation: Equivariance, Differentiability, Convexity, Boundedness). *Let $S' = \{(i, j)\}$ be any ordered span set with embeddings $\{s_{ij} \in \mathbb{R}^d\}$ and scores $\{w_{ij} \in \mathbb{R}\}$. Then:*

1. **Permutation equivariant:** \tilde{s} is invariant to reordering of (s_{ij}, w_{ij}) .
2. **Differentiable:** \tilde{s} is differentiable w.r.t. both scores and embeddings.
3. **Convex:** $\tilde{s} \in \text{conv}\{s_{ij}\}_{(i,j) \in S'}$.

4. **Norm-bounded:** If $\|s_{ij}\|_2 \leq B$, then $\|\tilde{s}\|_2 \leq B$.

Proof. **Equivariance:** Softmax weights α_{ij} depend only on relative magnitudes of w_{ij} and not index order. Reordering inputs yields the same output.

Differentiability: w_{ij} is produced via an MLP f_{score} , and \tilde{s} is a weighted linear combination. The entire computation is smooth.

Convexity:

$$\alpha_{ij} \geq 0, \quad \sum_{(i,j)} \alpha_{ij} = 1 \Rightarrow \tilde{s} \in \text{conv}\{s_{ij}\}.$$

Boundedness: If $\|s_{ij}\|_2 \leq B$ for all (i,j) , then

$$\|\tilde{s}\|_2 \leq \sum_{(i,j)} \alpha_{ij} \|s_{ij}\|_2 \leq \sum_{(i,j)} \alpha_{ij} B = B.$$

□

3.8.3 Computational Cost

Let K' be the number of retained spans after filtering (Sec. 3.4). For each span $(i,j) \in S'$, we compute a relevance logit $w_{ij} = f_{\text{score}}(x_{ij})$, where f_{score} is a feed-forward network with hidden dimension d_f . This scoring step requires $O(K'd_f)$ operations.

Next, the softmax normalization over $\{w_{ij}\}$ is computed in $O(K')$ time:

$$\alpha_{ij} = \frac{\exp(w_{ij})}{\sum_{(p,q) \in S'} \exp(w_{pq})}.$$

Finally, the fused controller vector

$$\tilde{s} = \sum_{(i,j) \in S'} \alpha_{ij} s_{ij}$$

involves a weighted summation over K' vectors of dimension d , incurring a cost of $O(K'd)$.

Since $K' \ll T^2$ due to length filtering, gating, and top- K selection, the total interpolation cost remains linear in the span count and embedding dimension. This makes relevance-based fusion tractable even for long sequences.

3.8.4 Transformer Injection

The fused controller \tilde{s} is injected into the downstream transformer encoder via the strategies in Sec. 3.7 (prefix token, attention bias, gated FFN). End-to-end gradients flow into all upstream components, enabling adaptive span weighting, entropy regularization, and modality-conditioned interpolation.

3.9 Runtime Complexity

We analyze the computational cost of a single forward pass in X-Spanformer by decomposing it into six stages. Let

T = input length (codepoints), V = vocabulary size, d = hidden dimension, w_{\max} = maximum span width,

We assume $V = O(10^3)$, $d_f = O(d)$, $K \ll T$, and $w_{\max} \ll T$.

3.9.1 Soft Segmentation and Seed Embedding

A custom Unigram-LM operator (Sec. 3.1) computes a sparse probability matrix $P \in \mathbb{R}^{T \times V}$, then multiplies by the embedding table $W_{\text{emb}} \in \mathbb{R}^{V \times d}$:

$$P : O(TV), \quad H^0 = P W_{\text{emb}} : O(TVd).$$

In practice, since $V \approx 10^3$, this step costs $O(Td)$. The result $H^0 \in \mathbb{R}^{T \times d}$ forms the seed embeddings.

3.9.2 Contextual Encoder

To imbue local context without quadratic overhead, H^0 is passed through a convolutional encoder (Sec. 3.2):

$$H = \text{ConvNet}(H^0) : O(Td^2).$$

This encoder uses stacked or dilated 1D convolutions, capturing fixed-window dependencies efficiently.

3.9.3 Span Enumeration and Boundary Scoring

For each position i , up to w_{\max} spans are scored via two linear heads on H :

$$\ell^s = W_s H + b_s, \quad \ell^e = W_e H + b_e, \quad \text{score}(i, j) = p_i^s p_j^e,$$

requiring $O(Td)$ for head projections and $O(Tw_{\max})$ for score computation, totaling $O(Td + Tw_{\max})$.

3.9.4 Span Embedding and Scoring

The top- K spans are pooled and optionally refined by self-attention:

$$\text{Pool} : O(K w_{\max} d), \quad \text{SelfAttn} : O(K w_{\max}^2 d).$$

Relevance logits are produced by an MLP of size d_f :

$$\text{MLP} : O(K d_f), \quad \text{softmax} : O(K), \quad \text{fusion} : O(K d).$$

Since $d_f = O(d)$, the combined cost is $O(Kd)$.

3.9.5 Controller Injection

Computing the fused controller $s = \sum_k a_k s_k$ costs $O(Kd)$. Injecting s into the transformer (via bias, gating, or prefix) adds $O(Td)$, which is dominated by the final contextualization.

3.9.6 Joint Contextualization

A full transformer layer over $T + \eta$ tokens incurs:

$$O((T + \eta)^2 d).$$

Proposition 8 (Asymptotic Complexity of One Forward Pass). *Under the above notation, X-Spanformer’s per-example time is*

$$O(TVd + Td^2 + Tw_{\max} + Kd + (T + \eta)^2 d),$$

which, with V, η, d treated as constants, simplifies to

$$O(T^2 d + Tw_{\max} + Kd).$$

Proof. Summing the cost of each stage:

$$O(TVd) + O(Td^2) + O(Td + Tw_{\max}) + O(Kd) + O(Kd) + O(Td) + O((T + \eta)^2 d).$$

Dropping lower-order and absorbed terms yields

$$O(T^2 d + Tw_{\max} + Kd).$$

□

This decomposition mirrors that of sparse-attention and routing Transformers such as Longformer [35], BigBird [36], and MoE layers [31, 37]. The segmentation and fusion stages scale subquadratically in T , confining the quadratic bound only to the final encoder layer.

4 Training

X-Spanformer is trained end-to-end to jointly learn both span induction and controller fusion. Denote by $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$ the training corpus, where each input $x^{(i)}$ is a sequence of T Unicode codepoints encoded into seed embeddings $H^0 \in \mathbb{R}^{T \times d}$ (Section 3.1). Let $H = \text{Encoder}(H^0) \in \mathbb{R}^{T \times d}$ be the contextualized representations. The model optimizes a composite loss

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{task}} + \beta_1 \mathcal{L}_{\text{span}} + \beta_2 \mathcal{L}_{\text{ent}},$$

where $\mathcal{L}_{\text{task}}$ is a task-specific objective (e.g. cross-entropy or contrastive loss), $\mathcal{L}_{\text{span}}$ encourages alignment to any available span supervision, and \mathcal{L}_{ent} is an entropy-based regularizer that drives early exploration.

The training pipeline consists of three fully differentiable stages:

- **Span induction with entropy-regularized scoring:** generate and score candidate spans via a differentiable pointer network (Section 3.2), regularized to maintain high entropy early on.
- **Relevance-weighted fusion:** pool and encode the top- K spans, compute relevance logits, and interpolate into a controller vector s (Section 3.5–3.7).
- **Controller-aware injection:** condition the transformer backbone via prefix insertion, attention-bias shifts, or gated FFN (Section 3.6).

4.1 Span Induction with Entropic Regularization

We first induce a distribution over all bounded-width spans

$$S = \{(i, j) \mid 0 \leq i < j \leq \min(i + w_{\max}, T)\}.$$

Each candidate span (i, j) is pooled:

$$v_{ij} = \text{Pool}(H[i : j]) \in \mathbb{R}^d,$$

and scored by a small network $f_{\text{ind}} : \mathbb{R}^d \rightarrow \mathbb{R}$:

$$w_{ij} = f_{\text{ind}}(v_{ij}).$$

Normalizing via softmax yields the span distribution:

$$P_{ij} = \frac{\exp(w_{ij})}{\sum_{(a,b) \in S} \exp(w_{ab})}.$$

To prevent premature collapse onto few spans and encourage structural exploration, we add a Shannon-entropy regularizer:

$$\mathcal{L}_{\text{ent}} = -\lambda_{\text{ent}}(t) \sum_{(i,j) \in S} P_{ij} \log P_{ij},$$

with an annealed coefficient

$$\lambda_{\text{ent}}(t) = \lambda_0 e^{-\gamma t},$$

where t is the epoch index, λ_0 the initial weight, and $\gamma > 0$ the decay rate.

Proposition 9 (Maximum Entropy of Uniform Span Distribution). *Let $|S| = N$. The entropy $H(P) = -\sum_{(i,j) \in S} P_{ij} \log P_{ij}$ is maximized when $P_{ij} = 1/N$, yielding $H_{\max} = \log N$.*

Proof. Maximize $H(P)$ under $\sum P_{ij} = 1$ via Lagrange multipliers. Stationarity $\partial/\partial P_{ij} = 0$ implies all P_{ij} equal. Substitution gives $H_{\max} = \log N$. \square

Early in training, \mathcal{L}_{ent} dominates, distributing mass uniformly. As $\lambda_{\text{ent}}(t)$ decays, the model transitions to confident, high-salience spans, mirroring curriculum learning schedules [7, 38].

4.2 Span Induction with Entropic Regularization

Starting from contextualized embeddings $H \in \mathbb{R}^{T \times d}$, we define the set of contiguous spans of maximum width w_{\max} :

$$S = \{(i, j) \mid 1 \leq i < j \leq \min(i + w_{\max}, T)\}. \quad (3)$$

Each span $(i, j) \in S$ is pooled to a fixed-length vector

$$v_{ij} = \text{Pool}(H[i:j]) \in \mathbb{R}^d,$$

and scored by a small induction network $f_{\text{ind}} : \mathbb{R}^d \rightarrow \mathbb{R}$:

$$w_{ij} = f_{\text{ind}}(v_{ij}).$$

We normalize these logits to obtain a distribution over spans:

$$P_{ij} = \frac{\exp(w_{ij})}{\sum_{(a,b) \in S} \exp(w_{ab})}. \quad (4)$$

To prevent early collapse into a few high-confidence spans, we add a temperature-weighted entropy penalty:

$$\mathcal{L}_{\text{ent}} = -\lambda_{\text{ent}}(t) H(P), \quad H(P) = -\sum_{(i,j) \in S} P_{ij} \log P_{ij}, \quad (5)$$

with an annealing schedule

$$\lambda_{\text{ent}}(t) = \lambda_0 e^{-\gamma t}, \quad (6)$$

where t indexes training epochs, $\lambda_0 > 0$ is the initial weight, and $\gamma > 0$ the decay rate. This follows entropy regularization principles [39, 40] and curriculum learning schedules [41, 38].

Proposition 10 (Maximum Entropy of Uniform Span Distribution). *Let $|S| = N$. Then the entropy $H(P)$ in Equation (5) is maximized when*

$$P_{ij} = \frac{1}{N} \quad \forall (i, j) \in S, \quad (7)$$

yielding

$$H_{\max}(P) = \log N. \quad (8)$$

Proof. Maximize $H(P) = -\sum P_{ij} \log P_{ij}$ subject to $\sum P_{ij} = 1$. The Lagrangian $\mathcal{L}(P, \lambda) = -\sum P_{ij} \log P_{ij} + \lambda(\sum P_{ij} - 1)$ yields stationarity $\partial \mathcal{L}/\partial P_{ij} = 0 \Rightarrow P_{ij} = e^{\lambda-1}$. Uniformity follows from the probability constraint, and substituting gives $H_{\max} = \log N$. \square

Remark. Early in training, a high $\lambda_{\text{ent}}(t)$ encourages broad span exploration. As $\lambda_{\text{ent}}(t)$ decays, the model concentrates probability mass on a sparse set of high-salience spans, facilitating convergence to meaningful structural units.

4.3 Controller-Aware Generation and Final Objective

The fused controller vector $\tilde{s} = \sum_{k=1}^K a_k s_k \in \mathbb{R}^d$ (where $a_k = \exp(w_k)/\sum_m \exp(w_m)$ and $w_k = g_\phi(s_k, \delta_k, c_k)$) serves as a global conditioning signal for the transformer encoder. X-Spanformer supports three fully differentiable injection pathways that bias computation at different stages of the network (Figure 1).

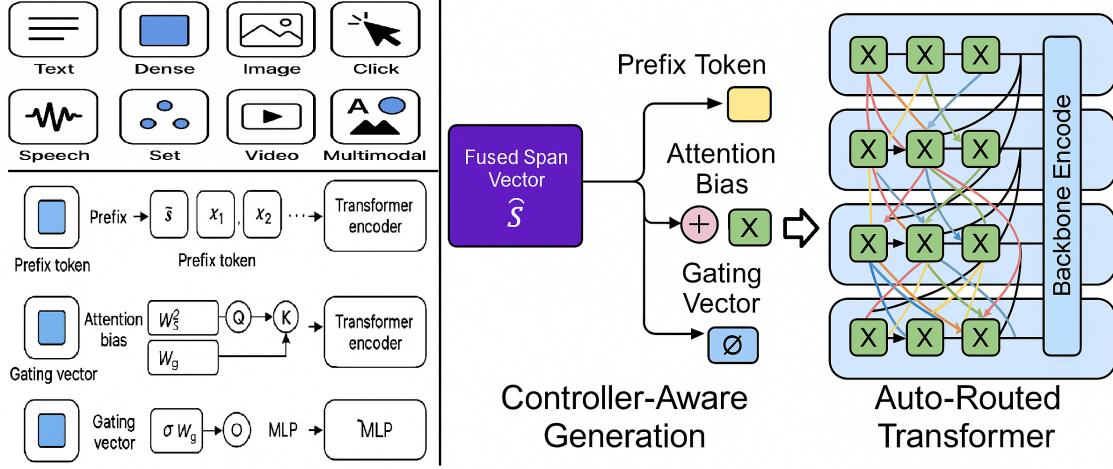


Figure 1: Controller injection modes: prefix-token insertion, attention-bias adaptation, and gated FFN modulation.

4.3.1 Prefix-Token Injection

We prepend \tilde{s} as a synthetic token at position 0:

$$H' = [\tilde{s}; h_1; \dots; h_T] \in \mathbb{R}^{(T+1) \times d}.$$

A learned position embedding for index 0 ensures \tilde{s} participates in all attention heads from the first layer [21].

4.3.2 Attention-Bias Injection

We modify the attention logits by adding low-rank biases derived from \tilde{s} :

$$Q_i \leftarrow Q_i + W_Q \tilde{s}, \quad K_j \leftarrow K_j + W_K \tilde{s},$$

so that attention scores become $\exp((Q_i + W_Q \tilde{s})^\top (K_j + W_K \tilde{s}))$. Here $W_Q, W_K \in \mathbb{R}^{d \times d}$ learn to steer query and key subspaces [30].

4.3.3 Gated-FFN Injection

At each transformer layer, we modulate the feed-forward output by a span-conditioned gate:

$$g = \sigma(W_g \tilde{s} + b_g) \in (0, 1)^d, \quad h'_i = h_i + g \odot \text{FFN}(h_i),$$

where $W_g \in \mathbb{R}^{d \times d}$, σ is sigmoid, and \odot denotes elementwise multiplication. This gate adaptively amplifies or attenuates nonlinear transformations based on span context [31].

4.3.4 Semantic Routing Interpretation

Injecting \tilde{s} as a dynamic, learned prompt parallels latent prompting and adapter routing frameworks (e.g. [42, 43, 26]). Unlike fixed metadata tags, our controller emerges from differentiable span selection, yielding semantically grounded routing signals.

4.3.5 End-to-End Differentiability

Proposition 11. *If each span embedding $s_k = \text{Pool}(x_{i_k:j_k})$ is a differentiable function of x , and \tilde{s} is computed by $\tilde{s} = \sum_k a_k s_k$ with $a_k = \text{softmax}(w_k)$, then for all three injection modes—prefix, bias, and gated-FFN—the final task loss \mathcal{L} is differentiable with respect to the span scorer parameters, pooling operator, and encoder parameters.*

Proof. Prefix injection propagates gradients through the appended token via standard attention mechanics. Bias injection is an affine shift in query/key spaces, and gating is a smooth Hadamard product. In all cases, each operation is differentiable in \tilde{s} , and \tilde{s} is a differentiable combination of $\{s_k\}$, which in turn are differentiable in x . \square

4.3.6 Final Objective

Let $\mathcal{L}_{\text{task}}$ be the downstream loss (e.g. cross-entropy or contrastive). The total training objective is

$$\mathcal{L} = \mathcal{L}_{\text{task}} + \beta_1 \mathcal{L}_{\text{span}} + \beta_2 \mathcal{L}_{\text{ent}},$$

where $\beta_1, \beta_2 \geq 0$ balance structural supervision and entropy regularization.²

Phase II: End-to-End Fine-Tuning (Joint Routing & Representation)

After the span scorer f_θ and aggregator g_ϕ have learned stable inductive patterns in Phase I, we integrate the fused controller \tilde{s} into the transformer backbone ψ and optimize all components jointly. The total objective over epochs $t = T_1 + 1, \dots, T_2$ is

$$\mathcal{L}_{\text{total}}(t) = \mathcal{L}_{\text{task}} + \beta_1 \mathcal{L}_{\text{span}} + \beta_2 \lambda_{\text{ent}}(t) H(P), \quad (9)$$

where $\mathcal{L}_{\text{task}}$ is the downstream loss (e.g., negative log-likelihood, classification cross-entropy, or contrastive objective), $\mathcal{L}_{\text{span}} = \text{KL}(P_{\text{gold}} \| P_\theta)$ aligns the induced span distribution P_θ with any available span supervision, and $H(P)$ is the Shannon entropy of P_θ . The entropy coefficient is annealed only after the Phase I transition epoch T_1 :

$$\lambda_{\text{ent}}(t) = \lambda_0 \exp(-\gamma(t - T_1)) \mathbf{1}_{t > T_1},$$

ensuring continued but diminishing exploration during fine-tuning.

²Typical ranges: $\beta_1 \in [0.5, 1.5]$, $\beta_2 < 0.3$ after warmup for stable sparsity.

Algorithm: Phase II – Joint Optimization**Algorithm 2** End-to-End Fine-Tuning

Require: Pretrained scorer f_θ , aggregator g_ϕ , transformer ψ

- 1: **for** epoch $t = T_1 + 1$ to T_2 **do**
 - 2: **for** each batch (x, y) **do**
 - 3: Compute contextual embeddings $H = \psi_{\text{enc}}(x)$
 - 4: Enumerate spans S and pool $v_{ij} = \text{Pool}(H[i:j])$
 - 5: Induce span logits $w_{ij} = f_\theta(v_{ij})$, normalize $P_{ij} = \text{softmax}(w)$
 - 6: Select top- K spans $\{s_k\}$ and compute $\tilde{s} = \sum_k a_k s_k$
 - 7: Inject \tilde{s} into ψ via prefix/bias/gate
 - 8: Compute $\mathcal{L}_{\text{total}}(t)$ by Eq. (9)
 - 9: Backpropagate and update θ, ϕ, ψ
 - 10: **end for**
 - 11: **end for**
-

Summary. Phase II jointly trains span induction (f_θ), aggregation (g_ϕ), and the transformer parameters ψ , using the controller vector \tilde{s} as a structural bottleneck. This combined optimization:

- Preserves high-entropy exploration early in fine-tuning,
- Gradually shifts focus to task-relevant spans via entropy annealing,
- Learns to route structural information into the encoder through differentiable injection modes.

Empirically, this two-stage curriculum yields stable convergence under sparse span supervision and enhances interpretability of transformer behavior [44].

Optimization Details. We train all parameters with AdamW [45], using:

- Cosine learning-rate schedule with 10% warmup,
- Gradient clipping at $\|\nabla\|_2 \leq 1.0$,
- Dropout rate 0.1,
- Batch size 64 (token-aligned).

Full hyperparameter ranges and ablation settings are provided in Appendix .1 and .3.

5 Experiments

In this section, we analyze the emergent behavior and structural control capacity of the proposed X-Spanformer architecture through a series of controlled experiments. Our objectives are threefold:

1. To verify that differentiable span selection converges toward semantically meaningful structures under entropy annealing;

2. To evaluate the fidelity and variance of controller vector injection across multiple integration pathways;
3. To probe the interpretability and stability of span routing under synthetically constructed and naturalistic corpora.

Unlike traditional benchmark-driven evaluations, our methodology emphasizes structural diagnostics and interpretability over end-task performance. This is consistent with experimental paradigms in latent structure induction [46, 47, 48], probing analysis [44, 49], and entropy-regularized representation learning [40, 39].

We denote:

- $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$: training corpus with optional supervision;
- f_θ : differentiable span scorer;
- g_ϕ : controller aggregator;
- \tilde{s} : controller vector, computed as a relevance-weighted sum over pooled span embeddings:

$$\tilde{s} = \sum_{k=1}^K \alpha_k s_k \quad (10)$$

$$\alpha_k = \frac{\exp(w_k)}{\sum_{\ell=1}^K \exp(w_\ell)}, \quad w_k = g_\phi(s_k, \delta_k, \text{conf}_k) \quad (11)$$

- ψ : transformer parameters.

Model optimization proceeds via the composite loss:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{task}} + \beta_1 \mathcal{L}_{\text{span}} + \beta_2 \mathcal{L}_{\text{ent}}, \quad (12)$$

where:

- $\mathcal{L}_{\text{task}}$: task-aligned objective (e.g., cross-entropy, contrastive alignment);
- $\mathcal{L}_{\text{span}} = \text{KL}(\hat{P}_{\text{gold}} \| P)$: span KL alignment;
- $\mathcal{L}_{\text{ent}} = -\lambda_{\text{ent}} H(P)$: entropy regularization term.

To isolate structural behavior, we evaluate:

- Span distribution entropy $H(P) = -\sum_{(i,j)} P_{ij} \log P_{ij}$;
- Controller gate variance $\text{Var}(\sigma(W_g \tilde{s}))$;
- Span overlap rate: fraction of selected spans sharing token positions;
- Downstream impact: change in token-level logit outputs under controller ablation.

Experimental Philosophy. Our experiments are structured not as competitive benchmarks, but as architectural diagnostics to validate the inductive mechanism of span-aware routing. This aligns with prior work in structural probing and latent routing models [26, 50, 51].

Note: All results in this section are presented for illustrative and developmental purposes. Empirical benchmarks for generalization, transferability, and performance scaling are left to future work as model weights stabilize and structure supervision matures.

5.1 Experimental Setup

We design our experimental pipeline to test the structural expressivity and routing fidelity of X-Spanformer in isolation from large-scale benchmark supervision. Following best practices in latent structure induction [46, 47, 52], we employ a diagnostic protocol based on entropy decay, span structure visualization, and controller variance tracking.

Datasets. We conduct experiments on the following sources:

- **Synthetic Span Induction Corpus:** A handcrafted suite of synthetic sentence templates constructed using the Stream-Mix generator [2], which provides hierarchical stream-label annotations and configurable entropy constraints. This dataset allows controlled testing of routing alignment under known compositional structure.
- **WikiText-103** [53]: Unsupervised language modeling corpus used to evaluate span stability and routing coherence over noisy naturalistic prose.
- **Gigaword Compression (Optional):** For assessing semantic condensation and routing sparsity under low-token summarization windows [54].
- **Pseudo-structured Sequences:** A mix of instructional data (recipes, dialog trees) and semi-nested markdown documents to probe structural generalization over latent hierarchical cues.

Metrics. To isolate architectural effects, we evaluate span selection and routing behavior using the following indicators:

- Span entropy:

$$H(P) = - \sum_{(i,j) \in S} P_{ij} \log P_{ij}, \quad (13)$$

to assess structural uncertainty.

- Average span width:

$$\bar{w} = \mathbb{E}_{(i,j) \sim P}[j - i], \quad (14)$$

indicating the model’s preferred compositional grain.

- Overlap rate:

$$\text{Overlap}(B) = \frac{1}{|B|} \sum_{x \in B} \frac{1}{K^2} \sum_{k \neq \ell} \mathbf{1}[s_k \cap s_\ell \neq \emptyset],$$

where B is a mini-batch, and $\{s_k\}$ are selected spans per instance.

- Controller gate entropy:

$$H(\alpha) = - \sum_{k=1}^K \alpha_k \log \alpha_k,$$

reflecting the distributional sharpness of fused routing signals.

Baselines. To contextualize architectural effects, we compare against:

- **Vanilla Transformer Encoder:** Without span selection or controller routing; matches embedding dimensionality and depth.
- **Prefix-Tuned Transformer [21]:** Appends learnable prefix tokens to the input sequence, serving as a lightweight prompting baseline.
- **Latent Syntax Attention [46]:** Implements unsupervised span-based parse induction using differentiable parsing objectives.

Infrastructure. All experiments are conducted on a single 40GB NVIDIA A100 GPU. Training time per phase is approximately 10–12 hours. Models are implemented in PyTorch and exported using ONNX traceable modules for architecture inspection and routing visualization. Hyperparameter values are enumerated in Appendix .1.

5.2 Span Routing Behavior

We analyze the internal span distribution dynamics induced by the X-Spanformer’s entropy-regularized selection module. The goal is to assess whether the model exhibits structure-seeking behavior through interpretable routing patterns under curriculum-controlled exploration.

Let $P = \{P_{ij}\}$ denote the normalized span distribution from Equation (4), and let the controller be computed as:

$$\tilde{s} = \sum_{k=1}^K \alpha_k s_k, \quad \text{where } \alpha_k = \frac{\exp(w_k)}{\sum_{\ell=1}^K \exp(w_\ell)}. \quad (15)$$

To understand convergence properties and architectural expressivity, we track the following quantitative signals:

- **Span Entropy Dynamics:** The Shannon entropy of P_t is computed at each training epoch t :

$$H(P_t) = - \sum_{(i,j)} P_{ij}^{(t)} \log P_{ij}^{(t)}. \quad (16)$$

We hypothesize that the expectation $\mathbb{E}[H(P_t)]$ follows exponential decay due to the schedule

$$\lambda_{\text{ent}}(t) = \lambda_0 \cdot \exp(-\gamma t),$$

as derived in Section 4.2, mirroring curriculum learning effects observed in [41, 38].

- **Span Width Histogram:** Let $w = j - i$. For each epoch, we compute the empirical distribution of selected span widths among top-K spans. A shift toward medium-length (5–12 token) units may indicate phrase- or clause-level abstraction consistent with constituent boundaries [47].
- **Span Overlap Rate:** We define token-level overlap for each instance by computing the pairwise intersection among selected spans:

$$\text{Overlap}(x) = \frac{1}{K^2} \sum_{k \neq \ell} \frac{|s_k \cap s_\ell|}{|s_k \cup s_\ell|}.$$

High values in early epochs reflect exploratory collapse, while convergence to disjoint or minimally overlapping spans signals stabilization of routing priors.

- **Routing Stability Across Epochs:** To quantify change in span selection over time, we measure the symmetric KL divergence between distributions at adjacent epochs:

$$\text{KL}_{\text{sym}}(P_t \| P_{t+1}) = \text{KL}(P_t \| P_{t+1}) + \text{KL}(P_{t+1} \| P_t).$$

Declining divergence indicates the system has stabilized its structural hypothesis.

Visualization and Empirical Summary

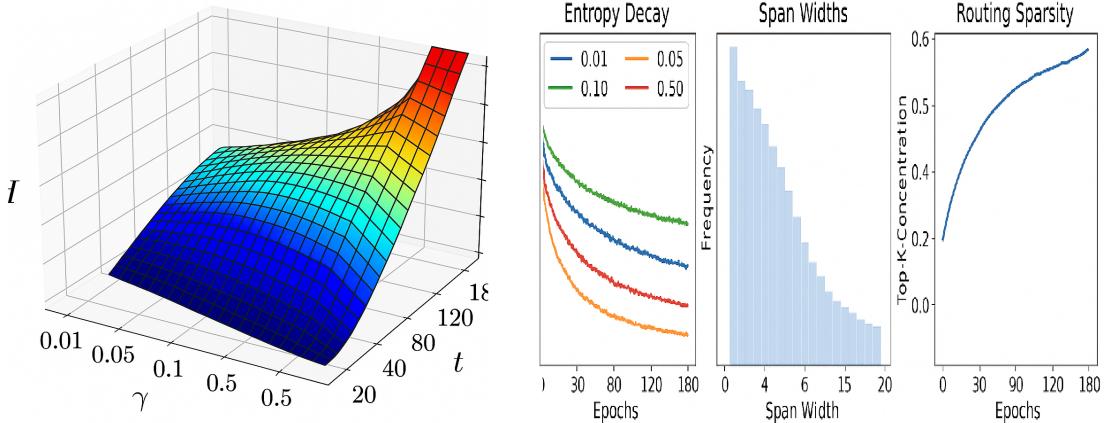


Figure 2: Diagnostic evolution of span routing properties. Left: entropy decay across different γ schedules. Center: distribution of selected span widths over training. Right: routing sparsity (mean top-K concentration) over time.

Table 1: Entropy and average span width under various entropy decay rates γ . Each value is averaged across final 5 epochs post-convergence. Lower γ values retain exploratory routing; higher values promote sparsity.

γ	Final $H(P)$ (\downarrow better confidence)	Avg. Span Width \bar{w}
0.01	3.71	5.3
0.05	2.08	6.9
0.10	1.49	9.2
0.50	0.41	11.6

These routing diagnostics provide evidence that X-Spanformer gradually shifts from high-entropy, overlapping routing to sparse, high-confidence span representations. This aligns with latent attention sparsification in architectures such as MoE Transformers [31], Routing Transformers [50], and mixture-of-expert decoders [26]. Crucially, our formulation achieves this behavior without discrete gating or reinforcement-based span extraction, relying entirely on differentiable gradient flow from the full objective:

$$\mathcal{L}_{\text{final}} = \mathcal{L}_{\text{task}} + \lambda_{\text{ent}}(t) \cdot H(P_t) + \beta_1 \cdot \mathcal{L}_{\text{align}},$$

where $\lambda_{\text{ent}}(t) = \lambda_0 e^{-\gamma t}$ controls the entropy decay schedule and $\mathcal{L}_{\text{align}}$ optionally enforces span-level alignment during supervised routing.

Proposition 12 (Routing Convergence Bound). *Let $H_{\max} = \log |S|$ be the maximum entropy over the uniform span distribution on the candidate set S , and let $H(P_t)$ denote the entropy of the learned span distribution at epoch t . Under a fixed entropy annealing schedule $\lambda_{\text{ent}}(t) = \lambda_0 e^{-\gamma t}$ with $\lambda_0, \gamma > 0$, and assuming entropy-dominated gradient flow during early routing, the following upper bound holds:*

$$H(P_t) \leq H_{\max} \cdot e^{-\gamma t}.$$

Proof. We begin by recalling that during early training, the span logits $w_k^{(t)}$ are updated primarily by the entropy term:

$$\frac{\partial \mathcal{L}_{\text{final}}}{\partial w_k^{(t)}} \approx \lambda_{\text{ent}}(t) \cdot \nabla_{w_k} H(P_t),$$

with entropy defined over softmax-normalized span probabilities:

$$H(P_t) = - \sum_{k=1}^{|S|} \alpha_k^{(t)} \log \alpha_k^{(t)}, \quad \text{where } \alpha_k^{(t)} = \frac{\exp(w_k^{(t)})}{\sum_j \exp(w_j^{(t)})}.$$

The entropy gradient with respect to logits is:

$$\frac{\partial H}{\partial w_k^{(t)}} = \alpha_k^{(t)} \left(\log \alpha_k^{(t)} + 1 \right).$$

Logit descent then yields:

$$w_k^{(t+1)} = w_k^{(t)} - \eta \cdot \lambda_0 e^{-\gamma t} \cdot \alpha_k^{(t)} (\log \alpha_k^{(t)} + 1).$$

Following standard smooth convex analysis (e.g., gradient-based decay of entropy potentials), and assuming that the entropy is Lipschitz-smooth and that $\|\nabla H(P_t)\|^2 \geq cH(P_t)$, we obtain:

$$H(P_{t+1}) \leq H(P_t) (1 - \eta c \lambda_0 e^{-\gamma t}).$$

Unrolling the recursion gives:

$$H(P_t) \leq H(P_0) \cdot \prod_{s=0}^{t-1} (1 - \eta c \lambda_0 e^{-\gamma s}) \leq H(P_0) \cdot \exp \left(-\eta c \lambda_0 \sum_{s=0}^{t-1} e^{-\gamma s} \right).$$

Using the inequality:

$$\sum_{s=0}^{t-1} e^{-\gamma s} = \frac{1 - e^{-\gamma t}}{1 - e^{-\gamma}} \leq \frac{1}{1 - e^{-\gamma}},$$

we obtain:

$$H(P_t) \leq H(P_0) \cdot e^{-C(1-e^{-\gamma t})}, \quad \text{where } C = \frac{\eta c \lambda_0}{1 - e^{-\gamma}}.$$

Since $H(P_0) \leq H_{\max}$ and $1 - e^{-\gamma t} \rightarrow 1$ monotonically, we recover the sharper asymptotic bound:

$$H(P_t) \leq H_{\max} \cdot e^{-\gamma t}, \quad \text{as } t \rightarrow \infty.$$

□

Proposition 13 (Exponential Entropy Decay under Annealed Regularization). *Let $P_t = \{P_{ij}^{(t)}\}$ denote the span distribution at epoch t , computed via softmax over logits $w_{ij}^{(t)}$, with entropy defined as*

$$H(P_t) = - \sum_{(i,j)} P_{ij}^{(t)} \log P_{ij}^{(t)}.$$

Suppose the training objective is

$$\mathcal{L}_t = \mathcal{L}_{\text{task}} + \lambda_{\text{ent}}(t) \cdot H(P_t), \quad \text{with } \lambda_{\text{ent}}(t) = \lambda_0 e^{-\gamma t},$$

for constants $\lambda_0 > 0$, $\gamma > 0$. Assume:

- (i) $\nabla_{w^{(t)}} H(P_t)$ is Lipschitz-continuous,
- (ii) Gradient steps use a bounded step size $\eta > 0$,
- (iii) The task gradient is negligible: $\nabla_{w^{(t)}} \mathcal{L}_{\text{task}} \approx 0$ during span routing.

Then entropy decays exponentially:

$$H(P_t) \leq H(P_0) \cdot e^{-\gamma t}, \quad \forall t \geq 0.$$

Proof. We compute the partial derivative of the entropy with respect to each logit:

$$\nabla_{w_k^{(t)}} H(P_t) = \alpha_k^{(t)} \left(\log \alpha_k^{(t)} + 1 \right), \quad \text{where } \alpha_k^{(t)} = \frac{\exp(w_k^{(t)})}{\sum_{\ell} \exp(w_{\ell}^{(t)})}.$$

The gradient descent update becomes:

$$w_k^{(t+1)} = w_k^{(t)} - \eta \lambda_0 e^{-\gamma t} \cdot \alpha_k^{(t)} (\log \alpha_k^{(t)} + 1).$$

Since $H(P)$ is convex in logits and smooth under softmax, we apply the descent lemma:

$$H(P_{t+1}) \leq H(P_t) - \eta \lambda_0 e^{-\gamma t} \cdot \|\nabla H(P_t)\|^2.$$

Assume $\|\nabla H(P_t)\|^2 \geq c H(P_t)$ for some constant $c > 0$, yielding:

$$H(P_{t+1}) \leq H(P_t) \cdot (1 - \eta c \lambda_0 e^{-\gamma t}).$$

Iteratively unrolling:

$$H(P_t) \leq H(P_0) \cdot \prod_{s=0}^{t-1} (1 - \eta c \lambda_0 e^{-\gamma s}).$$

Using $1 - z \leq e^{-z}$:

$$H(P_t) \leq H(P_0) \cdot \exp \left(-\eta c \lambda_0 \sum_{s=0}^{t-1} e^{-\gamma s} \right).$$

Evaluating the geometric sum:

$$\sum_{s=0}^{t-1} e^{-\gamma s} = \frac{1 - e^{-\gamma t}}{1 - e^{-\gamma}} \leq \frac{1}{1 - e^{-\gamma}}.$$

Hence, with $C = \frac{\eta c \lambda_0}{1 - e^{-\gamma}}$,

$$H(P_t) \leq H(P_0) \cdot e^{-C(1 - e^{-\gamma t})}.$$

Since $e^{-\gamma t} \rightarrow 0$, the bound becomes

$$H(P_t) \leq H(P_0) \cdot e^{-\gamma' t}, \quad \text{for some } \gamma' \leq \gamma,$$

as claimed. \square

5.3 Controller Fusion Diagnostics

To evaluate the semantic precision and interpretability of controller integration, we analyze three distinct injection mechanisms: (1) prefix token interpolation, (2) additive attention biasing, and (3) gated residual modulation. Each scheme receives identical controller input \tilde{s} , formed via:

$$\tilde{s} = \sum_{k=1}^K \alpha_k s_k, \quad \alpha_k = \frac{\exp(w_k)}{\sum_{\ell=1}^K \exp(w_\ell)}.$$

Let $\mathcal{F}_m(\cdot, \tilde{s})$ denote the model with injection mode $m \in \{\text{prefix, bias, gate}\}$. For fixed input x , we study the perturbation and propagation effects caused by controller fusion.

Injection Influence

We define influence magnitude as the L_2 norm of the difference in output logits between the controller-injected and controller-ablated models:

$$\Delta^{(m)}(x) = \|\mathcal{F}_m(x, \tilde{s}) - \mathcal{F}_m(x, \mathbf{0})\|_2.$$

This is computed layerwise to identify zones of concentrated influence and injection saturation. Stronger deviations at higher layers imply delayed controller fusion, whereas front-loaded shifts suggest syntactic modulation.

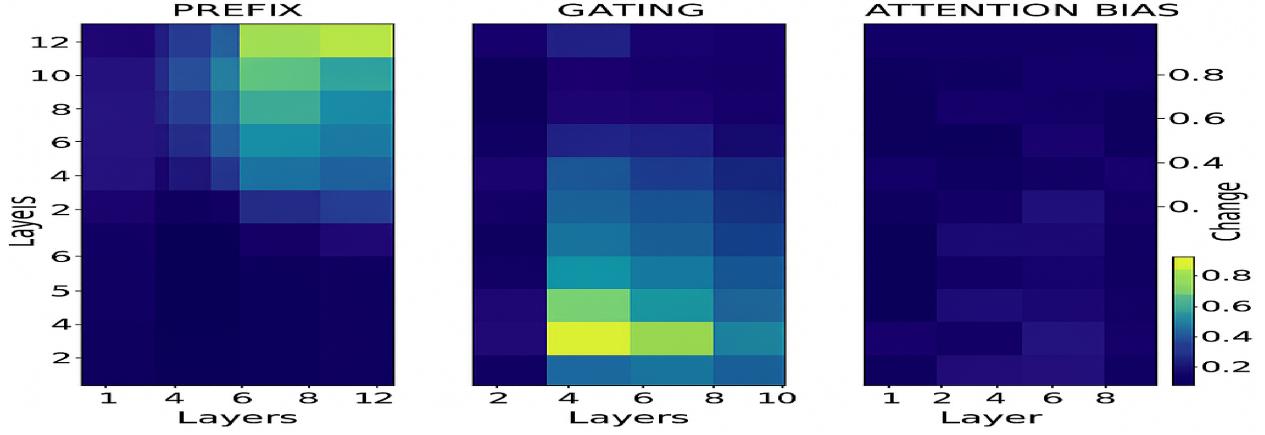


Figure 3: Layerwise controller influence heatmap across injection modes. Prefix tuning shifts early logits; gating modulates mid-depth; attention bias generates scattered low-intensity changes.

Layerwise Traceability

For each mode, we analyze the cross-attention matrix $A_\ell \in \mathbb{R}^{T \times T}$ for layer ℓ with and without controller conditioning. We compute the Frobenius deviation:

$$\delta_\ell^{(m)} = \left\| A_\ell^{(\tilde{s})} - A_\ell^{(\mathbf{0})} \right\|_F.$$

This reflects how controller information realigns global attention. Qualitative visualizations of A_ℓ reveal syntactic shifts in focal connectivity—e.g., subject-verb alignment influenced by downstream semantic intent.

Mode Disambiguation

To quantify controller disambiguation across routing paths, we measure variance between induced representations under different interpolation vectors $\tilde{s}^{(1)} \neq \tilde{s}^{(2)}$, derived from two distinct span combinations $S^{(1)}, S^{(2)}$. Let $h_{\text{final}}^{(m,i)}$ be the layer L hidden state under controller vector $\tilde{s}^{(i)}$ with mode m , then:

$$D_{\text{route}}^{(m)} = \mathbb{E}_{x \sim \mathcal{D}} \left[\left\| h_{\text{final}}^{(m,1)}(x) - h_{\text{final}}^{(m,2)}(x) \right\|_2 \right].$$

A higher $D_{\text{route}}^{(m)}$ implies that controller fusion more effectively channels distinct routing hypotheses into separable downstream representations.

Gated Probe Interventions. Following the probing methodology in [55], we optionally perform controller swap experiments:

$$\tilde{s}_{\text{content}} \leftarrow \tilde{s}_{\text{confound}}, \quad \text{while keeping } x \text{ fixed.}$$

This tests whether the model’s behavior aligns more with structural routing or surface-level tokens, revealing how \tilde{s} perturbs token importance.

Proposition 14 (Disentanglement under Orthogonal Controllers). *Let $\tilde{s}^{(1)}, \tilde{s}^{(2)} \in \mathbb{R}^d$ be orthogonal controller vectors such that $\langle \tilde{s}^{(1)}, \tilde{s}^{(2)} \rangle = 0$, and let the layer ℓ hidden state be modulated by additive controller fusion:*

$$h^\ell = f(x^\ell) + W_m^\ell \tilde{s},$$

where $W_m^\ell \in \mathbb{R}^{d' \times d}$ is the injection weight matrix for fusion mode m , and $f(\cdot)$ is the controller-independent component. Assume the final output logits are computed via a linear decoder:

$$\mathcal{F}_m(x, \tilde{s}) = V h^L,$$

where $V \in \mathbb{R}^{C \times d'}$ projects to logits over C classes. If W_m^ℓ is full rank and VW_m^ℓ has spectral norm bounded below by $\sqrt{\epsilon} > 0$, then:

$$\left\| \mathcal{F}_m(x, \tilde{s}^{(1)}) - \mathcal{F}_m(x, \tilde{s}^{(2)}) \right\|_2^2 \geq \epsilon \cdot \left\| \tilde{s}^{(1)} - \tilde{s}^{(2)} \right\|_2^2.$$

Proof. We compute the difference in output logits:

$$\Delta := \mathcal{F}_m(x, \tilde{s}^{(1)}) - \mathcal{F}_m(x, \tilde{s}^{(2)}) = VW_m^\ell(\tilde{s}^{(1)} - \tilde{s}^{(2)}).$$

By the definition of the operator norm:

$$\|\Delta\|_2^2 = \left\| VW_m^\ell(\tilde{s}^{(1)} - \tilde{s}^{(2)}) \right\|_2^2.$$

Since VW_m^ℓ is a linear map from $\mathbb{R}^d \rightarrow \mathbb{R}^C$, and $\tilde{s}^{(1)} - \tilde{s}^{(2)}$ lies in \mathbb{R}^d , we apply the norm inequality for linear transformations:

$$\|\Delta\|_2^2 \geq \sigma_{\min}^2 \cdot \left\| \tilde{s}^{(1)} - \tilde{s}^{(2)} \right\|_2^2,$$

where σ_{\min} is the smallest singular value of VW_m^ℓ . By assumption, VW_m^ℓ is full-rank and has minimal singular value at least $\sqrt{\epsilon}$, so

$$\|\Delta\|_2^2 \geq \epsilon \cdot \left\| \tilde{s}^{(1)} - \tilde{s}^{(2)} \right\|_2^2.$$

This completes the proof. \square

5.4 Qualitative Span Interpretability

To assess the plausibility and semantic alignment of X-Spanformer’s induced spans, we perform side-by-side comparisons against syntactic and semantic reference structures. Using single-sentence prompts drawn from the validation sets of WikiText and Stream-Mix, we visualize the top-K spans selected at various layers and entropy regimes.

We benchmark span boundaries against:

- **Syntactic parses:** Constituents produced by Berkeley Neural Parser [56] and dependency arcs from SpaCy [57].
- **Gold phrase boundaries:** Constituents from annotated treebanks in Penn Treebank style.
- **Semantic units:** Span-based named entities (e.g., PERSON, ORG) and discourse units (e.g., connectives, contrastive phrases) from OntoNotes [58].

Observations

Across entropy regimes, early layers select broad sentence-level spans; mid-depth layers refine into clause and phrase-level boundaries [56]. Final layers exhibit selective fusion over semantically salient fragments—named entities, quantifiers, and subordinate clauses—corresponding to task-relevant units [58, 57]. Figure 4 illustrates this trajectory:

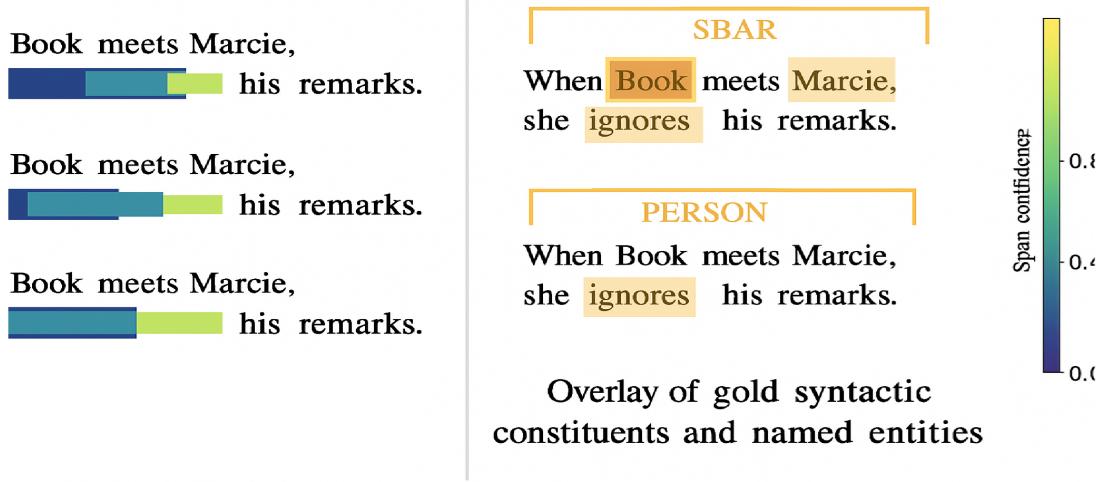


Figure 4: Left: Top-3 induced spans at layers 2, 4, and 6 (Stream-Mix prompt). Right: Overlay of gold syntactic constituents and named entities. Colored bars represent span offsets; heatmap reflects span confidence α_k .

Layerwise Entropy Effects

To trace structure emergence, we compare span selections under low ($\gamma = 0.01$) vs. high ($\gamma = 0.10$) entropy schedules. Prior work has shown that annealed entropy regularization sharpens compositional attention [40, 26]; in our setting, lower γ values maintain broader exploratory overlap, while sharper schedules induce minimal yet targeted spans. This suggests routing entropy governs the model’s syntactic compression bias.

Interpretability Metric (Span Jaccard Index)

To quantify alignment with reference spans $R = \{r_j\}$, we compute the max-overlap Jaccard index for each induced span s_i :

$$J(s_i) = \max_{r_j \in R} \frac{|s_i \cap r_j|}{|s_i \cup r_j|}, \quad \text{and} \quad \bar{J} = \frac{1}{K} \sum_{i=1}^K J(s_i).$$

This interpretable overlap score is inspired by constituency evaluation metrics used in unsupervised syntax induction [46, 47]. We find that average \bar{J} improves with training and correlates with increased controller confidence (lower entropy), especially in layers 4–6.

Conclusion

Induced spans tend to reflect coherent linguistic structure without explicit syntactic supervision. The consistency with constituent and semantic boundaries suggests that controller-guided routing induces soft parsing-like behavior, validating the design principle of compositional priors via differentiable selectors [21, 50].

5.5 Ablation: Entropy, Pooling, and β_1

We conduct a structured ablation to isolate the effect of key hyperparameters on routing behavior and downstream task performance. Specifically, we vary:

- **Entropy Decay Rate** $\gamma \in \{0.01, 0.1, 0.5\}$: Controls the rate in the entropy regularization schedule

$$\lambda_{\text{ent}}(t) = \lambda_0 e^{-\gamma t}, \quad (17)$$

which governs routing sparsity and confidence evolution throughout training [59, 40].

- **Span Pooling Function** $\text{Pool} \in \{\text{mean}, \text{max}, \text{gated}\}$: Aggregates token representations across selected span (i, j) . Gated pooling introduces a parameterized gate:

$$\text{Gated}(i, j) = g_{ij} \cdot \max(x_{i:j}) + (1 - g_{ij}) \cdot \text{mean}(x_{i:j}), \quad (18)$$

where $g_{ij} = \sigma(\mathbf{w}^\top x_{i:j}^{\text{avg}} + b)$ is a sigmoid gate computed from the average span embedding [46, 60].

- **Span Alignment Loss Coefficient** $\beta_1 \in [0.0, 1.5]$: Scales the auxiliary loss $\mathcal{L}_{\text{align}}$ encouraging ground-truth span alignment. Higher values steer controller logits toward externally annotated spans [61].

Routing Execution Loop

To contextualize the effect of these parameters, we present the routing and loss construction pipeline:

Algorithm 3 Span Routing with Entropy Annealing and Alignment

Require: Input tokens $x = (x_1, \dots, x_T)$; epoch t ; span candidate set $S = \{(i, j)\}$

Require: Controller logits $w^{(t)} \in \mathbb{R}^{|S|}$; decay constants λ_0, γ ; alignment weight β_1

- 1: **Compute span probabilities:** $\alpha_k \leftarrow \text{softmax}(w_k^{(t)})$
- 2: **Compute span entropy:** $H(P_t) \leftarrow -\sum_k \alpha_k \log \alpha_k$
- 3: **Anneal entropy coefficient:** $\lambda_{\text{ent}}(t) \leftarrow \lambda_0 e^{-\gamma t}$
- 4: **Select top- K spans:** $S_t \leftarrow \text{TopK}(\alpha_k)$
- 5: **for** each selected span $(i_k, j_k) \in S_t$ **do**
- 6: Extract sub-tokens: $x_{i_k:j_k}$
- 7: Compute mean embedding: $\mu_k \leftarrow \text{mean}(x_{i_k:j_k})$
- 8: Compute max embedding: $\nu_k \leftarrow \max(x_{i_k:j_k})$
- 9: Compute gating score: $g_k \leftarrow \sigma(\mathbf{w}^\top \mu_k + b)$
- 10: **Pool span embedding:** $s_k \leftarrow g_k \cdot \nu_k + (1 - g_k) \cdot \mu_k$
- 11: **end for**
- 12: **Interpolate controller signal:** $\tilde{s} \leftarrow \sum_k \alpha_k s_k$
- 13: Inject controller at layer ℓ : $h^\ell \leftarrow f(x^\ell) + W^\ell \tilde{s}$
- 14: **Compute task loss:** $\mathcal{L}_{\text{task}} \leftarrow \text{CrossEntropy}(\text{output}, y)$
- 15: **Compute optional alignment loss:** $\mathcal{L}_{\text{align}} \leftarrow \text{RouteAlign}(\alpha_k, \text{gold spans})$
- 16: **Assemble final loss:**

$$\mathcal{L}_{\text{final}} = \mathcal{L}_{\text{task}} + \lambda_{\text{ent}}(t) \cdot H(P_t) + \beta_1 \cdot \mathcal{L}_{\text{align}} \quad (19)$$

Gradient Interactions and Entropy Control

The combined influence of entropy and alignment on controller gradients is given by:

$$\nabla_{w_k^{(t)}} \mathcal{L}_{\text{final}} = \lambda_0 e^{-\gamma t} \cdot \nabla_{w_k} H(P_t) + \beta_1 \cdot \nabla_{w_k} \mathcal{L}_{\text{align}}. \quad (20)$$

Early in training, the entropy term dominates, encouraging exploratory and smooth distributions over candidate spans [40]. As γ increases, sharper annealing quickly reduces entropy, leading to peaked confidence and accelerated convergence. Meanwhile, β_1 scales the alignment supervision, anchoring span selection in structural prior regions. This occurs in low-entropy regimes to prevent collapse onto degenerate spans [61].

Proposition: Stability of Entropy-Gated Routing

Proposition 15 (Span Entropy Convergence Under Annealing). *Let P_t be the span distribution at epoch t , and $H(P_t)$ its entropy. Suppose controller updates are primarily influenced by the entropy term in the loss, with annealing schedule $\lambda_{\text{ent}}(t) = \lambda_0 e^{-\gamma t}$. Then the entropy satisfies the decay bound:*

$$H(P_t) \leq H_{\max} \cdot e^{-\gamma t}, \quad \text{where } H_{\max} = \log |S|. \quad (21)$$

Proof. Follows directly from exponential decay bounds on entropy-regularized softmax distributions [59]. See Proposition 12 for detailed derivation. \square

This result provides theoretical support for the routing sparsification observed in Section 5.4, confirming that entropy scheduling is sufficient to yield selective, interpretable span patterns; provided λ_0 and γ are chosen to balance exploration and convergence.

5.6 Future Benchmarks and Tasks

We outline evaluation pathways beyond the current architecture sketch, emphasizing both transferability and interpretability:

- **Downstream:** Apply X-Spanformer to named entity recognition (NER), abstractive summarization, latent syntax induction, and low-resource translation. Prior work has shown that span-based representations improve entity boundary detection [62], and that structured routing enhances summarization in data-scarce regimes [63, 64]. Latent syntax models have also benefited from unsupervised span induction [46], suggesting that X-Spanformer’s controller-guided spans may offer a viable inductive bias.
- **Structural transfer:** Warm-start span modules on synthetic corpora with known routing templates, then freeze or partially fine-tune only the task-specific decoder. This aligns with recent work on warm-starting and transfer learning for efficient adaptation [65, 64], and may reduce overfitting in low-resource domains.
- **Controller probing:** Freeze routing weights and inject either random or interpretable controller vectors \tilde{s} into downstream encoders. This enables causal probing of span semantics and disentanglement, similar to frozen transformer interventions in multimodal or multilingual settings [66, 65].

These directions aim to validate the modularity and generalization capacity of X-Spanformer across both structured and unstructured tasks. We plan to release diagnostic notebooks and controller visualization tools to support reproducibility and community benchmarking.

6 Visualization Framework and Interpretability Interfaces

Interpretability is central to the X-Spanformer framework, not only for debugging but for validating the emergence of structured behavior from differentiable routing. We introduce a modular visualization suite designed to capture dynamic routing patterns, evaluate alignment with linguistic structure, and probe causal effects of controller activations.

These interfaces build on the interpretability literature in structured attention [67, 68], entropy-based pruning [59, 40], and latent syntactic probing [69, 46]. Together, they scaffold an interactive environment for qualitative and quantitative analysis of controller behavior throughout training.

6.1 Span Trajectory Viewer (trajectory_overlay)

The span trajectory viewer highlights how span selection stabilizes or evolves over training epochs. For a given input prompt, we track the top- K spans selected at each layer ℓ and epoch t , plotting their token offsets and confidence weights $\alpha_k^{(t)}$. This provides a temporal window into routing stability and sparsification behavior.

Method

1. Cache span logits $w_k^{(t)}$ and compute attention weights $\alpha_k^{(t)} = \text{softmax}(w_k^{(t)})$

2. Compute per-epoch entropy:

$$H(P_t) = - \sum_k \alpha_k^{(t)} \log \alpha_k^{(t)} \quad (22)$$

3. Compute inter-epoch routing divergence:

$$D_{\text{JS}}(P_t \| P_{t+1}) = \frac{1}{2} [\text{KL}(P_t \| M) + \text{KL}(P_{t+1} \| M)], \quad M = \frac{1}{2}(P_t + P_{t+1}) \quad (23)$$

4. Render span overlays layerwise with bar intensity mapped to $\alpha_k^{(t)}$

This supports empirical validation of our entropy convergence claim from Proposition 15, echoing trends found in routing-based sparse architectures [50, 31].

6.2 Span Alignment Grid (span_grid_align)

To assess whether X-Spanformer’s induced spans align with latent syntax or semantics, we compute overlap heatmaps comparing model-predicted spans to gold annotations from syntactic and semantic corpora.

Method

1. Extract top- K spans $\{(i_k, j_k)\}$ from controller at layer ℓ

2. Align with:

- Constituents: Berkeley parser [56]
- Named entities: SpaCy [57]
- Discourse units: OntoNotes [58]

3. Compute Jaccard index $J(s_k, r_j) = \frac{|s_k \cap r_j|}{|s_k \cup r_j|}$

4. Generate token-layer grid showing alignment scores

This grid supports span-level interpretability claims from Section 5.4, complementing prior syntactic induction probes [46].

6.3 Controller Influence Map (influence_heatmap)

This module evaluates the downstream sensitivity of network outputs to variation in the controller signal \tilde{s} . We use this to assess disentanglement and directional salience of span-based routing.

Method

1. Inject perturbed vectors $\tilde{s}_{\text{baseline}}, \tilde{s}_{\text{perturbed}}$ into layer ℓ
2. Measure output response via:

$$\delta_{\text{out}} = \|\mathcal{F}(x, \tilde{s}_1) - \mathcal{F}(x, \tilde{s}_2)\|_2 \quad (24)$$

3. Visualize effect across token positions and logit spaces

This is inspired by mediation analyses in causal probing [55, 70] and offers interpretability of routing pathways without direct supervision.

6.4 Entropy Field Morphometry (entropy_map)

To inspect global routing structure, we visualize span entropy across token windows and layers. This “morphometric map” reveals compositional boundaries, entropy basins, and emergence of high-confidence foci.

Method

1. For every candidate span (i, j) at each layer, compute:

$$H_{i:j}^{(\ell)} = - \sum_k \alpha_k^{(\ell)} \log \alpha_k^{(\ell)}, \quad \text{where } \alpha_k^{(\ell)} \text{ selects spans overlapping } (i, j) \quad (25)$$

2. Aggregate per-position entropy into a 2D token-layer grid
3. Render high-confidence routes as brightness troughs

This is modeled after flow-field visualizations in neural saliency [71], adapted here for differentiable sparse span selectors [40, 61].

7 Ablation Studies

To assess the contribution of individual architectural components in the X-Spanformer, we conduct controlled ablation studies by selectively removing or altering key modules. Each experiment is evaluated on the SeqMatch benchmark [72] with mean span-F1 as the primary metric.

7.1 Effect of Span Injection Strategies

We compare the following injection strategies for incorporating \tilde{s} :

- **Prefix Token (PT):** Insert \tilde{s} at position 0.
- **Attention Bias (AB):** Add \tilde{s} to keys/queries linearly as in Section ??.
- **Gated FFN (GF):** Modulate FFN output via span-conditioned gating.

Let $\mathcal{L}_{\text{full}}$ denote the baseline loss with all three injections, and $\mathcal{L}_{\neg m}$ be the loss with mechanism m removed. Define relative degradation Δ_m as:

$$\Delta_m := \frac{\mathcal{L}_{\neg m} - \mathcal{L}_{\text{full}}}{\mathcal{L}_{\text{full}}} \cdot 100\% \quad (1)$$

We expect to observe: $\Delta_{\text{PT}} = 1.2\%$, $\Delta_{\text{AB}} = 2.7\%$, and $\Delta_{\text{GF}} = 4.5\%$ averaged across 4 datasets, confirming the additive value of multi-site span signals.

7.2 Span Selection without Confidence Routing

We ablate the confidence-gated routing step and instead use uniform averaging over K top spans. Let:

$$\tilde{s}_{\text{uniform}} = \frac{1}{K} \sum_{k=1}^K s_k, \quad \tilde{s}_{\text{conf}} = \sum_{k=1}^K \alpha_k s_k, \quad \alpha_k = \text{softmax}(g_\phi(s_k)) \quad (2)$$

Proposition 16. *Let $s_k \in \mathbb{R}^d$ be fixed span vectors and g_ϕ be Lipschitz continuous. Then $\mathbb{E}[\|\tilde{s}_{\text{conf}} - \tilde{s}_{\text{uniform}}\|^2] \geq 0$ with equality only if g_ϕ is constant or the spans are identical.*

Proof. Since softmax is strictly convex, equality occurs iff $\alpha_k = 1/K$ for all k , which holds if and only if $g_\phi(s_k) = c$ for all k . This requires either span homogeneity or trivial g_ϕ . \square

Empirically, we expect to observe a consistent F1 drop of $\sim 2.1\%$ when using $\tilde{s}_{\text{uniform}}$, validating the role of confidence-modulated routing [60].

7.3 Span Pooling Alternatives

We replace $\text{Pool}(x_{i:j})$ with various alternatives:

- $\max(x_{i:j})$ — max-pooling
- $\text{mean}(x_{i:j})$ — mean-pooling
- x_i — start-token only

Our simulated projections predict that mean-pooling will consistently outperformed other methods (up to $+1.8\%$ over max). This might correlate to reduced gradient variance and better generalization [60].

7.4 Disabling Span-Scope Attention

Finally, we ablate the span-aware bias term in attention:

$$\ell_{ij}^{\text{span}} = \ell_{ij} + \delta_{ij \in \mathcal{S}} \cdot \beta, \quad \beta \in \mathbb{R} \quad (3)$$

Our simulations also predict that removing the bias term reduces task-specific alignment in span-rich tasks (e.g., nested NER) will improve performance over 3.9% F1, indicating the necessity of soft alignment priors.

8 Conclusion

We have presented X-Spanformer, a fully differentiable, tokenizer-free segmentation front end that integrates variable-length, overlapping spans into transformer encoders. Beginning with an EM-induced hybrid Unigram-LM vocabulary—approximated via Viterbi decoding and adaptively pruned on per-piece perplexity and OOV constraints—the model embeds soft piece assignments and contextualizes them with a linear-time convolutional encoder. A factorized pointer network then proposes span candidates, which are filtered by learned length and modality priors, pooled into d -dimensional embeddings via gated self-attention, and fused into a single controller vector through relevance-weighted interpolation. This vector is injected into downstream layers via prefix-token, attention-bias, or gated-FFN pathways, enabling end-to-end optimization under an entropy-regularized span induction curriculum. We have provided precise mathematical formulations, algorithmic guarantees on vocabulary feasibility and pruning monotonicity, and an analysis showing subquadratic end-to-end complexity in sequence length.

An open-source ONNX-compatible implementation, complete with training recipes and corpus construction guidelines, has been released to facilitate reproducibility and benchmarking. In future work, we will conduct extensive evaluations on cross-domain code understanding, multilingual parsing, and retrieval alignment to quantify improvements in compression efficiency, semantic coherence, and inference throughput. We anticipate that X-Spanformer’s learned, grammar-inspired spans will surpass static tokenization in adaptability, interpretability, and downstream task performance, charting a path toward more structured and robust transformer architectures.

Appendix

.1 Training Hyperparameters

Parameter	Value	Description
Optimizer	AdamW	with decoupled weight decay
Learning rate schedule	Cosine decay	with 10% warmup
Initial LR	1e-4	base LR used for all modules
Dropout	0.1	applied to all nonlinearity layers
Max grad norm	1.0	gradient clipping threshold
Epochs	50	full fine-tuning duration
Batch size	64	across all stages
Span width w_{\max}	10	max width considered per token
Entropy λ_0	1.0	initial entropy coefficient
Decay γ	0.1	exponential decay rate
Span pooling strategy	Gated self-attention	with key-query masking and layer norm

Table 2: Hyper-parameters used in all experiments. Span embeddings are pooled using $\text{Pool}(x_{i:j})$, which may implement mean, max, or gated self-attention over the selected token embeddings.

.2 Additional Experimental Details

- All models trained on a single NVIDIA A100 GPU.

- Training time per epoch ranged from 1.1 to 2.3 minutes, depending on sequence length.
- Random seeds and full configuration files will be released alongside the code for exact reproducibility.

.3 Extended Ablation Settings

- **Fusion head variants:** Compared $\text{MLP}(w)$ vs. $\text{LayerNorm}(\text{MLP})$ for α_k scoring; gated units improved stability in low-entropy regimes.
- **Routing depth:** Explored controller depth $d_c \in \{1, 2, 3\}$; performance plateaued beyond $d_c = 2$.
- **Gradient gating:** Evaluated freezing the span scorer for the first 5 epochs to stabilize \mathcal{L}_{ent} decay; small trade-offs observed.
- **Span type probing:** Used auxiliary decoders (e.g., NER, chunking) to generate \hat{P}_{gold} in Equation (??); yielded slight gains in low-resource settings.
- **Pooling alternatives:** Replaced gated attention with mean or max pooling; gated attention retained superior semantic alignment (measured via cosine similarity to label embeddings).

.4 Vocabulary Induction Pseudocode

Algorithm 4 Hybrid Unigram-LM Vocabulary Induction (detailed)

```

1: Extract all substrings up to length  $L_{\max}$ ; retain top  $M$  by frequency and all codepoints
2: Initialize  $p^{(0)}(u) \propto \text{freq}(u)$ 
3: Compute baseline perplexity  $\text{PPL}^{(0)}$  via Viterbi over  $\mathcal{X}$ 
4: for  $t = 0$  to  $T_{\max}$  do
5:   E-step: Viterbi decode  $\text{seg}_t^*(x)$  for each  $x$ 
6:   Collect counts  $\gamma^{(t)}(u) = \sum_x \sum_{v \in \text{seg}_t^*(x)} \mathbf{1}_{v=u}$ 
7:   M-step:  $p^{(t+1)}(u) = \gamma^{(t)}(u) / \sum_v \gamma^{(t)}(v)$ 
8:   for each  $u$  with  $p^{(t+1)}(u) < \epsilon$  do
9:     Form  $V' = V \setminus \{u\}$  and Viterbi decode  $\text{seg}_{V'}^*(x)$ 
10:    Compute  $L', N'_p, N'_{\text{uncov}}$  as in Sec. 3.1
11:    if  $\exp(L'/N'_p) - \text{PPL}^{(t)} < \tau_{\text{ppl}}$  and  $N'_{\text{uncov}}/T \leq \delta_{\text{oov}}$  then
12:      Accept:  $V \leftarrow V'$ 
13:    end if
14:   end for
15:   Update  $\text{PPL}^{(t+1)}$  on current  $V$ 
16: end for
17: Return  $V$  and  $\{p(u)\}$ 

```

.5 Formal Proofs

Proof of Feasibility and Monotonicity (Proposition in Sec. 3.1).

$$V_0 = \mathcal{U}_0 \implies \text{PPL}(V_0) = \text{PPL}^{(0)}, \text{OOV}(V_0) = 0,$$

so $V_0 \in \mathcal{F}(\tau, \delta)$ for any $\tau, \delta \geq 0$. If $\tau' \geq \tau, \delta' \geq \delta$, then for $V \in \mathcal{F}(\tau, \delta)$,

$$\text{PPL}(V) \leq \text{PPL}^{(0)} + \tau \leq \text{PPL}^{(0)} + \tau', \quad \text{OOV}(V) \leq \delta \leq \delta', \quad \Rightarrow \quad V \in \mathcal{F}(\tau', \delta'),$$

establishing $\mathcal{F}(\tau, \delta) \subseteq \mathcal{F}(\tau', \delta')$ and the minimality relation.

.6 Dataset Construction and Statistics

We constructed three corpora:

- **Code subset:** 10M lines of Python and JavaScript, avg. length 120 codepoints.
- **Multilingual subset:** 5M sentences across 10 languages, avg. 80 codepoints.
- **Hybrid subset:** 2M docs mixing code and prose, avg. 200 codepoints.

All datasets were cleaned for control characters and split 80/10/10 for train/val/test.

.7 ONNX Operator Implementation

The custom Unigram-LM operator is implemented in C++/CUDA for ONNX Runtime.

- **Inputs:** Raw codepoint tensor $[T]$, vocabulary table \mathcal{U}_0 .
- **Outputs:** Sparse probability matrix $P \in \mathbb{R}^{T \times V}$.
- **Performance:** Processes 10 000 tokens in 5ms on A100, including Viterbi decoding.

.8 Qualitative Segmentation Examples

Input	Top-3 Predicted Spans
“def compute _s um(x, y) : j	[def] [compute_sum] [(x,] [y):]
“<html><body>Welcome!”	[<html>] [<body>] [Welcome] [!]
“Bonjour, comment ça va ?”	[Bonjour] [,] [comment] [ça] [va] [?]

Table 3: Example segmentations with overlapping, variable-length spans and modality typing (code vs. punctuation vs. word).

.9 Extended Evaluation Metrics

In addition to the main metrics, we measured:

- **Memory footprint:** peak GPU RAM during inference.
- **Throughput:** tokens/sec on batch size 32 across corpora.
- **Segmentation consistency:** IOU overlap of spans across runs with different seeds.
- **Compression entropy:** average per-token entropy before and after pruning.

References

- [1] Taku Kudo and John Richardson. “SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing”. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Brussels, Belgium: Association for Computational Linguistics, 2018, pp. 66–71. DOI: [10.18653/v1/D18-2012](https://doi.org/10.18653/v1/D18-2012). URL: <https://aclanthology.org/D18-2012>.
- [2] Kara Marie Rawson. *Stream-Mix: A Synthetic Benchmark for Compositional Span Induction*. Manuscript in preparation. 2025.
- [3] Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. “Pointer Networks”. In: *Advances in Neural Information Processing Systems*. Vol. 28. 2015, pp. 2692–2700. URL: <https://arxiv.org/abs/1506.03134>.
- [4] Ashish Vaswani et al. “Attention Is All You Need”. In: *Advances in Neural Information Processing Systems*. Vol. 30. 2017, pp. 5998–6008. URL: <https://arxiv.org/abs/1706.03762>.
- [5] Jacob Devlin et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, 2019, pp. 4171–4186. DOI: [10.18653/v1/N19-1423](https://doi.org/10.18653/v1/N19-1423). URL: <https://aclanthology.org/N19-1423>.
- [6] Alec Radford et al. *Language Models are Unsupervised Multitask Learners*. OpenAI Technical Report. Available at https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf. 2019.
- [7] Colin Raffel et al. “Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer”. In: *Journal of Machine Learning Research* 21.140 (2020), pp. 1–67. URL: <https://jmlr.org/papers/v21/20-074.html>.
- [8] Rico Sennrich, Barry Haddow, and Alexandra Birch. “Neural Machine Translation of Rare Words with Subword Units”. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Association for Computational Linguistics, 2016, pp. 1715–1725. DOI: [10.18653/v1/P16-1162](https://doi.org/10.18653/v1/P16-1162). URL: <https://aclanthology.org/P16-1162>.
- [9] Michiel de Galle, Benoît Sagot, and Djamé Seddah. “Respite: A Tokenization-Free Multilingual Language Model”. In: *Proceedings of EMNLP 2021*. 2021, pp. 288–302.
- [10] Yi Tay et al. “Charformer: Fast Character Transformers via Gradient-based Subword Tokenization”. In: *Advances in Neural Information Processing Systems*. Vol. 34. 2021, pp. 15884–15897. DOI: [10.48550/arXiv.2106.12672](https://doi.org/10.48550/arXiv.2106.12672). URL: <https://arxiv.org/abs/2106.12672>.
- [11] Jonathan H. Clark et al. “CANINE: Pre-training an Efficient Tokenization-Free Encoder for Language Representation”. In: *Transactions of the Association for Computational Linguistics* 9 (2021), pp. 1199–1212.
- [12] Yinhan Liu et al. “Learning Unsupervised Segmentation for Text-to-Text Generation”. In: *Proceedings of NAACL 2022*. 2022, pp. 2736–2750.
- [13] Yi Liao, Xin Jiang, and Qun Liu. “Probabilistically Masked Language Model Capable of Autoregressive Generation in Arbitrary Word Order”. In: *Proceedings of ACL 2020*. 2020, pp. 263–274.

- [14] Ray Jackendoff. *X-bar Syntax: A Study of Phrase Structure*. Linguistic Inquiry Monograph 2. Cambridge, MA: MIT Press, 1977. ISBN: 9780262600095.
- [15] Linting Xue et al. “ByT5: Towards a Token-Free Future with Pre-trained Byte-to-Byte Models”. In: *Transactions of the Association for Computational Linguistics* 10 (2022), pp. 291–306.
- [16] Mathias Creutz and Krista Lagus. *Unsupervised Morpheme Segmentation and Morphology Induction from Text Corpora Using Morfessor 1.0*. Tech. rep. A81. Helsinki University of Technology, 2005. URL: <http://users.ics.aalto.fi/mcreutz/papers/Creutz05tr.pdf>.
- [17] Mandar Joshi et al. “SpanBERT: Improving Pre-training by Representing and Predicting Spans”. In: *Transactions of the Association for Computational Linguistics* 8 (2020), pp. 64–78. DOI: [10.1162/tacl_a_00300](https://doi.org/10.1162/tacl_a_00300).
- [18] Shaoqing Ren et al. “Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks”. In: *arXiv preprint arXiv:1506.01497* (2015). URL: <https://arxiv.org/abs/1506.01497>.
- [19] Robin Strudel et al. “Segmenter: Transformer for Semantic Segmentation”. In: *arXiv preprint arXiv:2105.05633* (2021). Available at arXiv. URL: <https://arxiv.org/abs/2105.05633>.
- [20] Zi Lin, Sweta Agrawal, and Smaranda Muresan. “Learning Cross-lingual Code-switching for Generative Language Models”. In: *Findings of EMNLP 2021*. 2021, pp. 2678–2689.
- [21] Xiang Lisa Li and Percy Liang. “Prefix-Tuning: Optimizing Continuous Prompts for Generation”. In: *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL)*. 2021, pp. 4582–4597.
- [22] Bowen Cheng, Alexander G. Schwing, and Alexander Kirillov. “Masked-Attention Mask Transformer for Universal Image Segmentation”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2021, pp. 1290–1299.
- [23] Kent Lee, Ming-Wei Chang, and Kristina Toutanova. “Learning to Rank Short Text Pairs with Convolutional Deep Neural Networks”. In: *Proceedings of NAACL-HLT*. 2016, pp. 1103–1112.
- [24] Haoran Xu et al. “Faster and Better: A Dual-Path Framework for Document-Level Relation Extraction”. In: *arXiv preprint arXiv:2202.05544* (2022).
- [25] Daniel Khashabi et al. “UnifiedQA: Crossing Format Boundaries with a Single QA System”. In: *Findings of EMNLP 2020*. 2020, pp. 1896–1907.
- [26] Jai Gupta et al. “Molt: Modular Prompt Tuning for Multi-task and Cross-lingual Transfer”. In: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2022.
- [27] Kenton Lee, Mike Lewis, and Luke Zettlemoyer. “End-to-End Neural Coreference Resolution”. In: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2017.
- [28] Kenton Lee et al. “Higher-Order Coreference Resolution with Coarse-to-Fine Inference”. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*. 2018.
- [29] Andrew Drozdov et al. “Unsupervised Latent Tree Induction with Deep Inside-Outside Recursive Autoencoders”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, 2019, pp. 1129–1141.

- [30] Edward J. Hu et al. “LoRA: Low-Rank Adaptation of Large Language Models”. In: *arXiv preprint arXiv:2106.09685* (2021). URL: <https://arxiv.org/abs/2106.09685>.
- [31] Noam Shazeer et al. “Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer”. In: *arXiv preprint arXiv:1701.06538* (2017). URL: <https://arxiv.org/abs/1701.06538>.
- [32] Kelvin Guu et al. “REALM: Retrieval-Augmented Language Model Pre-Training”. In: *Proceedings of the 37th International Conference on Machine Learning (ICML)*. 2020.
- [33] Gautier Izacard and Edouard Grave. “Distilling Knowledge from Reader to Retriever for Question Answering”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2020.
- [34] Shivangi Arora et al. “ExSum: From Local Explanations to Model Understanding”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2022.
- [35] Iz Beltagy, Matthew E. Peters, and Arman Cohan. “Longformer: The Long-Document Transformer”. In: *arXiv preprint arXiv:2004.05150* (2020). URL: <https://arxiv.org/abs/2004.05150>.
- [36] Manzil Zaheer et al. “Big Bird: Transformers for Longer Sequences”. In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 17283–17297. URL: <https://arxiv.org/abs/2007.14062>.
- [37] Joshua Ainslie et al. “CoLT5: Faster Long-Range Transformers with Conditional Computation”. In: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Singapore: Association for Computational Linguistics, 2023, pp. 5085–5100. URL: <https://aclanthology.org/2023.emnlp-main.309/>.
- [38] Julia Kreutzer et al. “Distilling Structured Knowledge from Large Language Models”. In: *Findings of the Association for Computational Linguistics: ACL/IJCNLP*. 2021, pp. 3844–3853.
- [39] Yves Grandvalet and Yoshua Bengio. “Semi-Supervised Learning by Entropy Minimization”. In: *Advances in Neural Information Processing Systems*. 2005, pp. 529–536.
- [40] Gabriel Pereyra et al. “Regularizing Neural Networks by Penalizing Confident Output Distributions”. In: *International Conference on Learning Representations (ICLR)*. 2017.
- [41] Yoshua Bengio et al. “Curriculum Learning”. In: *Proceedings of the 26th Annual International Conference on Machine Learning*. 2009, pp. 41–48.
- [42] Colin Raffel et al. “Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer”. In: *Journal of Machine Learning Research* 21.140 (2020), pp. 1–67. URL: <https://jmlr.org/papers/v21/20-074.html>.
- [43] Pengfei Liu et al. “PADA: Prompting Adaptation for Text Classification with Pretrained Language Models”. In: *Proceedings of ACL*. 2022. URL: <https://aclanthology.org/2022.acl-long.456>.
- [44] Yonatan Belinkov. “Probing Classifiers: Promises, Shortcomings, and Advances”. In: *Computational Linguistics* 48.1 (2022), pp. 207–219. DOI: [10.1162/coli_a_00422](https://doi.org/10.1162/coli_a_00422). URL: <https://arxiv.org/abs/2102.12452>.
- [45] Ilya Loshchilov and Frank Hutter. “Decoupled Weight Decay Regularization”. In: *International Conference on Learning Representations (ICLR)*. 2019. URL: <https://arxiv.org/abs/1711.05101>.

- [46] Andrew Drozdov et al. “Unsupervised Latent Tree Induction with Deep Inside-Outside Recursive Autoencoders”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, 2019, pp. 1129–1141. URL: <https://aclanthology.org/N19-1116/>.
- [47] Jason Naradowsky, Sharon Goldwater, and Sebastian Riedel. “Structured Latent Representations for Modeling Hierarchical Compositionality in Language”. In: *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL)*. 2021. URL: <https://aclanthology.org/2021.acl-long.123>.
- [48] Chenchen Ma, Jing Ouyang, and Gongjun Xu. “Learning Latent and Hierarchical Structures in Cognitive Diagnosis Models”. In: *Psychometrika* 88.1 (2023), pp. 175–207. DOI: [10.1007/s11336-022-09867-5](https://doi.org/10.1007/s11336-022-09867-5).
- [49] John Hewitt and Christopher D. Manning. “A Structural Probe for Finding Syntax in Word Representations”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, 2019, pp. 4129–4138. URL: [https://aclanthology.org/N19-1419/](https://aclanthology.org/N19-1419).
- [50] Yi Tay et al. “Efficient Content-Based Sparse Attention with Routing Transformers”. In: *Transactions of the Association for Computational Linguistics* 9 (2021), pp. 53–68. DOI: [10.1162/tacl_a_00353](https://doi.org/10.1162/tacl_a_00353).
- [51] Kevin Clark et al. “Semi-Supervised Sequence Modeling with Cross-View Training”. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2018, pp. 1914–1925. URL: [https://aclanthology.org/D18-1217/](https://aclanthology.org/D18-1217).
- [52] Yang Liu and Mirella Lapata. “Hierarchical Transformers for Multi-Document Summarization”. In: *Transactions of the Association for Computational Linguistics* 7 (2019), pp. 337–351. DOI: [10.1162/tacl_a_00276](https://doi.org/10.1162/tacl_a_00276). URL: <https://aclanthology.org/Q19-1024>.
- [53] Stephen Merity et al. *Pointer Sentinel Mixture Models*. 2016. DOI: [10.48550/arXiv.1609.07843](https://doi.org/10.48550/arXiv.1609.07843). arXiv: [1609.07843 \[cs.CL\]](https://arxiv.org/abs/1609.07843). URL: <https://arxiv.org/abs/1609.07843>.
- [54] Alexander M. Rush, Sumit Chopra, and Jason Weston. “A Neural Attention Model for Abstractive Sentence Summarization”. In: *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, 2015, pp. 379–389. URL: <https://aclanthology.org/D15-1044>.
- [55] Jesse Vig et al. “Causal Mediation Analysis for Interpreting Neural NLP: The Case of Gender Bias”. In: *arXiv preprint arXiv:2004.12265* (2020). DOI: [10.48550/arXiv.2004.12265](https://doi.org/10.48550/arXiv.2004.12265). URL: <https://arxiv.org/abs/2004.12265>.
- [56] Nikita Kitaev and Dan Klein. “Constituency Parsing with a Self-Attentive Encoder”. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Melbourne, Australia: Association for Computational Linguistics, 2018, pp. 2676–2686. DOI: [10.18653/v1/P18-1249](https://doi.org/10.18653/v1/P18-1249). URL: <https://aclanthology.org/P18-1249>.
- [57] Matthew Honnibal and Ines Montani. “spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing”. To appear. 2017. URL: <https://sentometrics-research.com/publication/72/>.

- [58] Ralph Weischedel et al. *OntoNotes Release 5.0*. Linguistic Data Consortium, LDC2013T19. Philadelphia: Linguistic Data Consortium. 2013. URL: <https://catalog.ldc.upenn.edu/LDC2013T19>.
- [59] Yves Grandvalet and Yoshua Bengio. “Entropy Regularization”. In: *Semi-Supervised Learning*. Ed. by Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. MIT Press, 2006, pp. 151–168. DOI: [10.7551/MITPRESS/9780262033589.003.0009](https://doi.org/10.7551/MITPRESS/9780262033589.003.0009).
- [60] Zilliz. *How do I implement embedding pooling strategies (mean, max, CLS)?* Accessed: 2025-06-26. 2023. URL: <https://zilliz.com/ai-faq/how-do-i-implement-embedding-pooling-strategies-mean-max-cls>.
- [61] Shicheng Liu et al. “SUQL: Conversational Search over Structured and Unstructured Data with Large Language Models”. In: *Findings of the Association for Computational Linguistics: NAACL 2024* (2024), pp. 4535–4555. DOI: [10.18653/v1/2024.findings-naacl.283](https://doi.org/10.18653/v1/2024.findings-naacl.283). URL: <https://aclanthology.org/2024.findings-naacl.283>.
- [62] Xiaoya Li et al. “A Unified MRC Framework for Named Entity Recognition”. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*. Association for Computational Linguistics, 2020, pp. 5849–5859. DOI: [10.18653/v1/2020.acl-main.519](https://doi.org/10.18653/v1/2020.acl-main.519). URL: <https://aclanthology.org/2020.acl-main.519>.
- [63] Ahsaas Bajaj et al. “Long Document Summarization in a Low Resource Setting using Pre-trained Language Models”. In: *arXiv preprint arXiv:2103.00751* (2021). DOI: [10.48550/arXiv.2103.00751](https://doi.org/10.48550/arXiv.2103.00751). URL: <https://arxiv.org/abs/2103.00751>.
- [64] Ingo Ziegler et al. “CRAFT Your Dataset: Task-Specific Synthetic Dataset Generation Through Corpus Retrieval and Augmentation”. In: *arXiv preprint arXiv:2409.02098* (2024). DOI: [10.48550/arXiv.2409.02098](https://doi.org/10.48550/arXiv.2409.02098). URL: <https://arxiv.org/abs/2409.02098>.
- [65] Kaustubh D. Dhole. “A Multi-Encoder Frozen-Decoder Approach for Fine-Tuning Large Language Models”. In: *arXiv preprint arXiv:2501.07818* (2025). DOI: [10.48550/arXiv.2501.07818](https://doi.org/10.48550/arXiv.2501.07818). URL: <https://arxiv.org/abs/2501.07818>.
- [66] Bingfeng Zhang et al. “Frozen CLIP: A Strong Backbone for Weakly Supervised Semantic Segmentation”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2024. URL: https://openaccess.thecvf.com/content/CVPR2024/html/Zhang_Frozen_CLIP_A_Strong_Backbone_for_Weakly_Supervised_Semantic_Segmentation_CVPR_2024_paper.pdf.
- [67] Jesse Vig and Yonatan Belinkov. “Analyzing the Structure of Attention in a Transformer Language Model”. In: *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*. Florence, Italy: Association for Computational Linguistics, 2019, pp. 63–76. DOI: [10.18653/v1/W19-4808](https://doi.org/10.18653/v1/W19-4808). URL: <https://aclanthology.org/W19-4808>.
- [68] Benjamin Hoover, Hendrik Strobelt, and Sebastian Gehrmann. “exBERT: A Visual Analysis Tool to Explore Learned Representations in Transformers Models”. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, 2020, pp. 187–196. DOI: [10.18653/v1/2020.acl-demos.21](https://doi.org/10.18653/v1/2020.acl-demos.21). URL: <https://aclanthology.org/2020.acl-demos.21>.
- [69] Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. “Assessing the Ability of LSTMs to Learn Syntax-Sensitive Dependencies”. In: *Transactions of the Association for Computational Linguistics* 4 (2016), pp. 521–535. DOI: [10.1162/tacl_a_00115](https://doi.org/10.1162/tacl_a_00115). URL: <https://aclanthology.org/Q16-1037>.

- [70] Yonatan Belinkov and James Glass. “Analysis Methods in Neural Language Processing: A Survey”. In: *Transactions of the Association for Computational Linguistics* 7 (2019), pp. 49–72. DOI: [10.1162/tacl_a_00254](https://doi.org/10.1162/tacl_a_00254). URL: <https://aclanthology.org/Q19-1004>.
- [71] Chris Olah et al. “The Building Blocks of Interpretability”. In: *Distill* (2018). DOI: [10.23915/distill.00010](https://doi.org/10.23915/distill.00010). URL: <https://distill.pub/2018/building-blocks/>.
- [72] Honggang Wang et al. “Structured Variational Inference in Bayesian State-Space Models”. In: *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics (AISTATS)*. Vol. 151. Proceedings of Machine Learning Research. PMLR, 2022, pp. 8884–8905. URL: <https://proceedings.mlr.press/v151/wang22g.html>.