# Causal Decoding

Complete Decoder-Only Language Model: Embeddings, Transformer Blocks,
Output Projection, and Autoregressive Generation with Sampling

The culmination of Phase 3: a complete CausalLM that takes token IDs
and generates token IDs, the same interface as any production LLM.
This is the naive version -- full forward pass recompute at every
generation step -- deliberately inefficient to motivate KV caching.

This demo covers:
1. Full forward pass walkthrough with shape tracing
2. Causal property verification (THE centerpiece)
3. Sampling strategies: greedy, temperature, top-k, top-p
4. Autoregressive generation loop visualization
5. Computational cost: naive $O(n^2)$ vs KV cache $O(n)$
6. Model parameter analysis for Llama/Mistral configs

Model config: V=256, d=64, layers=2, h=4, h_kv=2, d_ff=172
Random seed: 42
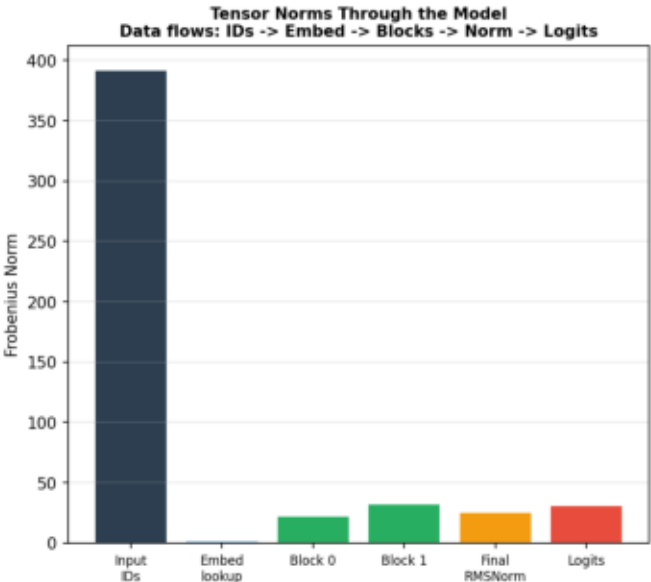Number of visualizations: 6

Examples: 6

# Summary of Findings

1. Forward Pass: Token IDs (B,L) -> embedding lookup (B,L,d) ->
   N transformer blocks -> final RMSNorm -> output logits (B,L,V).
   Weight tying: W_out = E.T shares memory, saving V*d parameters.
   Softmax of last-position logits gives next-token probabilities.

2. Causal Property (CENTERPIECE): For sequences sharing prefix [0..k-1],
   logits at positions 0..k-1 are EXACTLY identical regardless of tokens
   at positions k+. Verified: changing token at position 3 has ZERO effect
   on positions 0-2 (diff < 1e-12). This enables autoregressive generation.

3. Sampling Strategies: Temperature controls sharpness (T<1: deterministic,
   T>1: random). Top-k keeps only k highest logits. Top-p (nucleus) keeps
   the smallest set with cumulative probability >= p. Combined pipeline:
   logits -> /T -> top-k -> top-p -> softmax -> categorical sample.

4. Generation Loop: Naive version does full forward pass at each step,
   processing the growing sequence (P, P+1, ..., P+n-1 tokens).
   Total cost: n*P + n(n-1)/2 token-steps. With KV cache: P + n.
   The redundant recomputation motivates KV caching.

5. Computational Cost: Naive FLOPs grow super-linearly with generation
   length. KV cache eliminates redundant K/V projections for previous
   tokens. Speedup ratio grows approximately linearly with n.
   This is the single most impactful optimization in LLM inference.

6. Parameter Analysis: Transformer blocks are >86% of total parameters.
   Llama 2 7B: ~6.74B. Llama 3 8B: ~8.03B. Mistral 7B: ~7.24B.
   None use weight tying; tying would save V*d (131M-525M).
   FFN dominates per-block params (~67% MHA, ~81% GQA).

# Example 1: Full Forward Pass Walkthrough

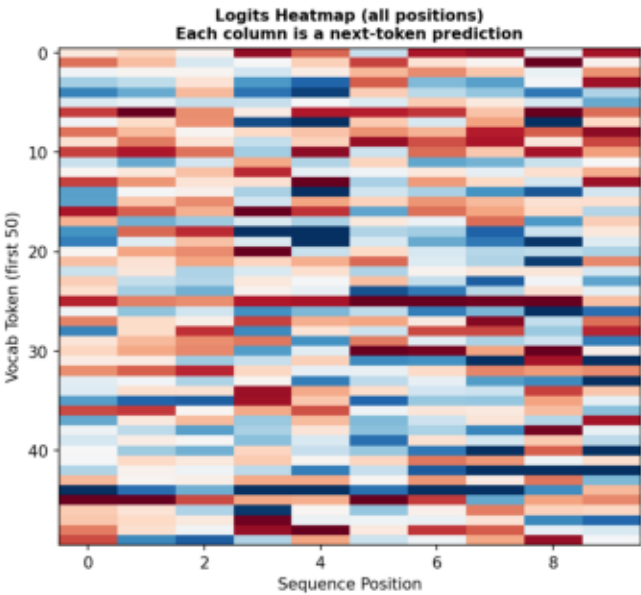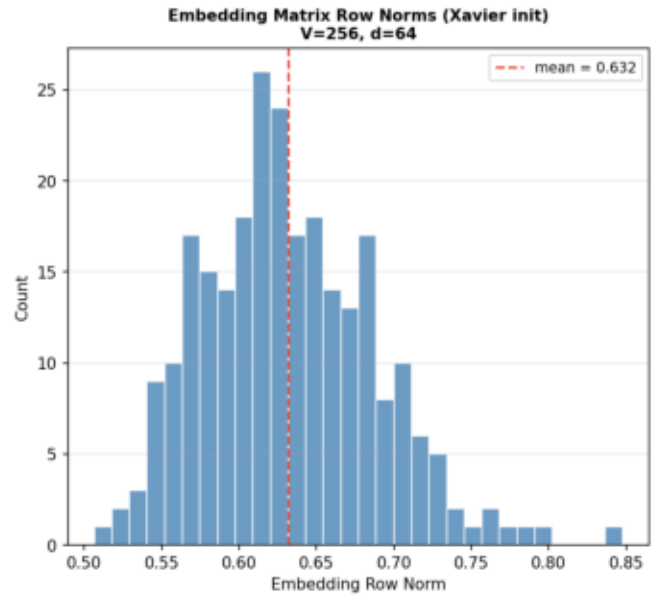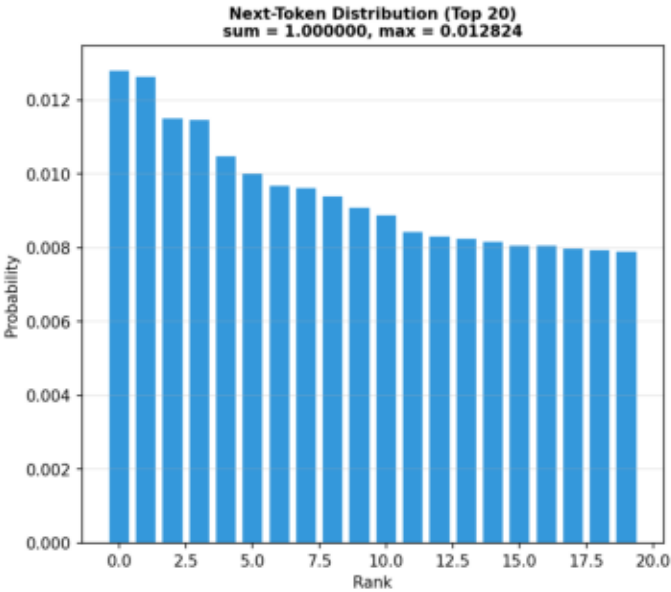

**Causal Decoding: Full Forward Pass Walkthrough**

**Tensor Norms Through the Model**
Data flows: IDs -> Embed -> Blocks -> Norm -> Logits

**Next-Token Distribution (Top 20)**
sum = 1.000000, max = 0.012824

```
COMPLETE SHAPE TABLE
====================

Token IDs:       (1, 10)
Embedding:       (1, 10, 64)
After block i:   (1, 10, 64)
Final RMSNorm:   (1, 10, 64)
Logits:          (1, 10, 256)

Next-token probs:
  logits[:, -1, :] -> softmax
  Shape: (1, 256)

Weight tying:
  E.shape = (256, 64)
  W_out = E.T = (64, 256)
  logits = x_norm @ E.T
  Shares memory: True
```
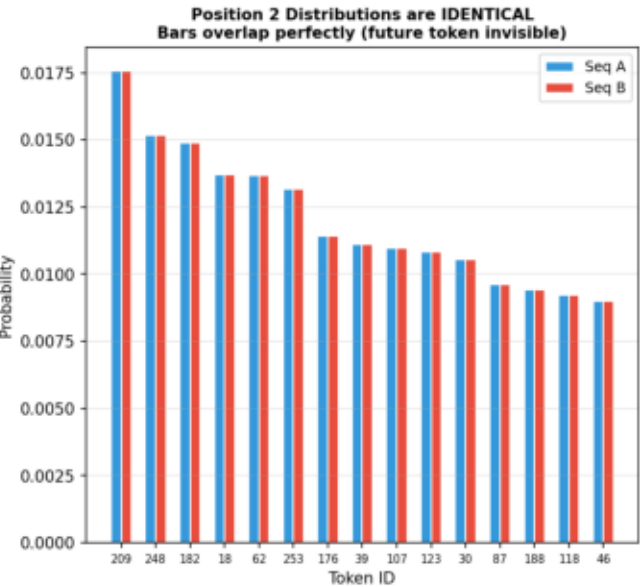
**Embedding Matrix Row Norms (Xavier init)**
V=256, d=64

mean = 0.632

**Logits Heatmap (all positions)**
Each column is a next-token prediction

```
ARCHITECTURE SUMMARY
====================

CausalLM(
  embedding:  E[256, 64]
  blocks:     2 x TransformerBlock(
                d=64, h=4,
                h_kv=2, d_ff=172
              )
  final_norm: RMSNorm(64)
  W_out:      E.T (weight tying)
)

Forward: E[tokens] -> blocks -> norm -> logits
Generate: forward -> sample -> append -> repeat

CAVEAT: Random weights produce
random (meaningless) distributions.
Trained models learn meaningful
next-token predictions.
```
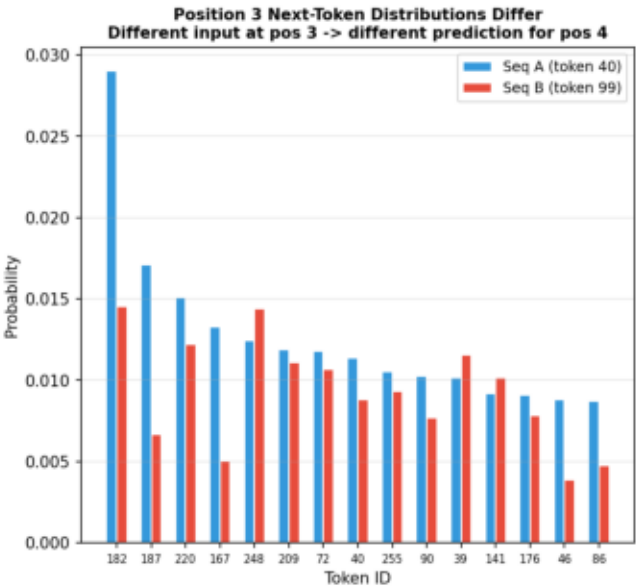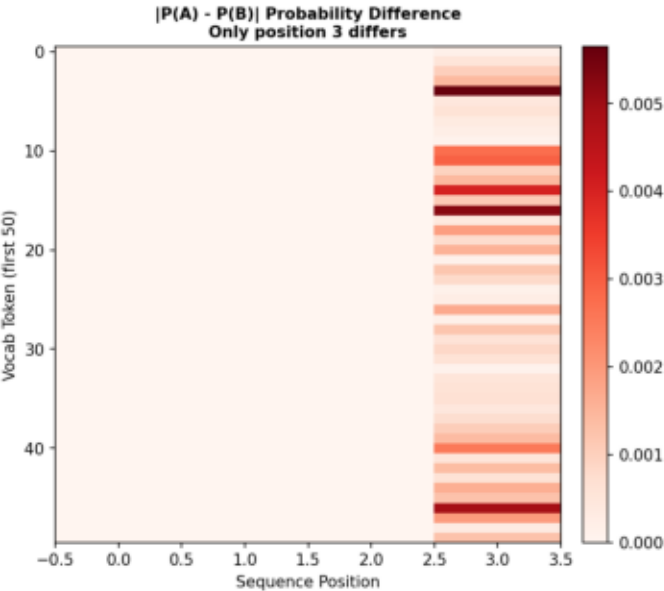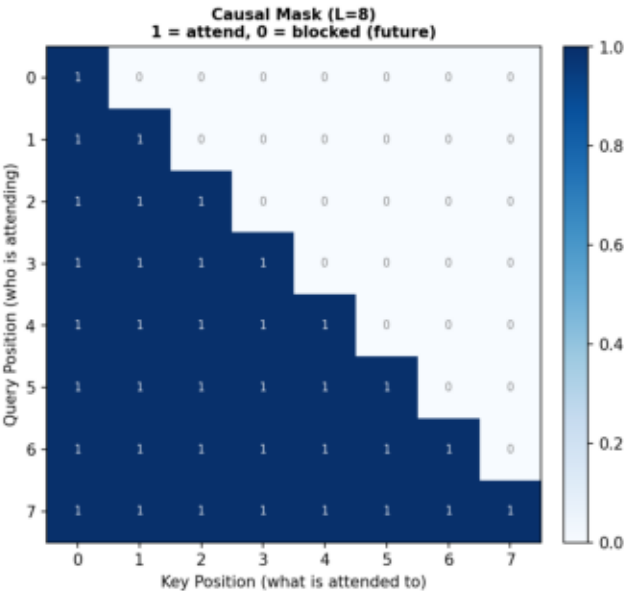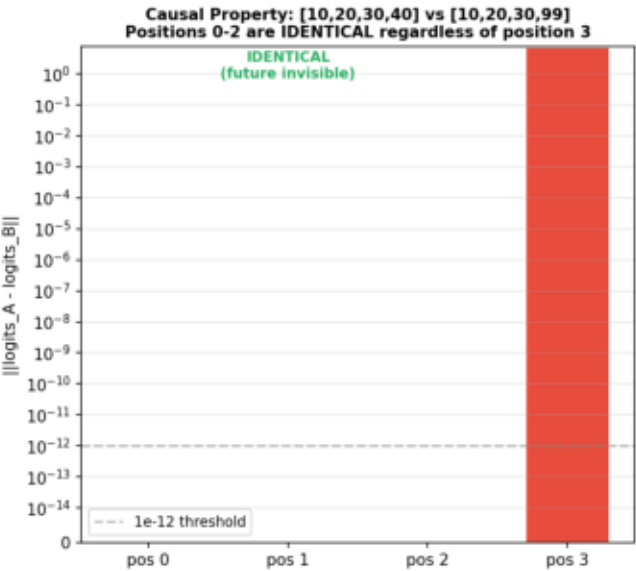
# Example 2: Causal Property Verification (CENTERPIECE)

## Causal Property Verification: Future Tokens Are Invisible



**Causal Property: [10,20,30,40] vs [10,20,30,99]**
Positions 0-2 are IDENTICAL regardless of position 3

**Causal Mask (L=8)**
1 = attend, 0 = blocked (future)

**|P(A) - P(B)| Probability Difference**
Only position 3 differs

**Position 3 Next-Token Distributions Differ**
Different input at pos 3 -> different prediction for pos 4

**Position 2 Distributions are IDENTICAL**
Bars overlap perfectly (future token invisible)

```
THE CAUSAL PROPERTY
===================

For sequences sharing prefix [0..k-1]:
   logits[0..k-1] are IDENTICAL
   regardless of tokens at positions k+.

WHY: The causal attention mask zeros
out all future positions. Position i
can only attend to positions 0..i.

CONSEQUENCE: We can generate left-
to-right, one token at a time.
Each new token depends only on the
prefix, never on future tokens.

This enables autoregressive generation:
   1. Forward pass on prompt
   2. Sample next token from last logits
   3. Append token, repeat

VERIFIED:
   Prefix positions: max diff = 0.0e+00
   Divergent position: diff = 7.1171
   Single vs full seq: diff = 6.6e-15
```
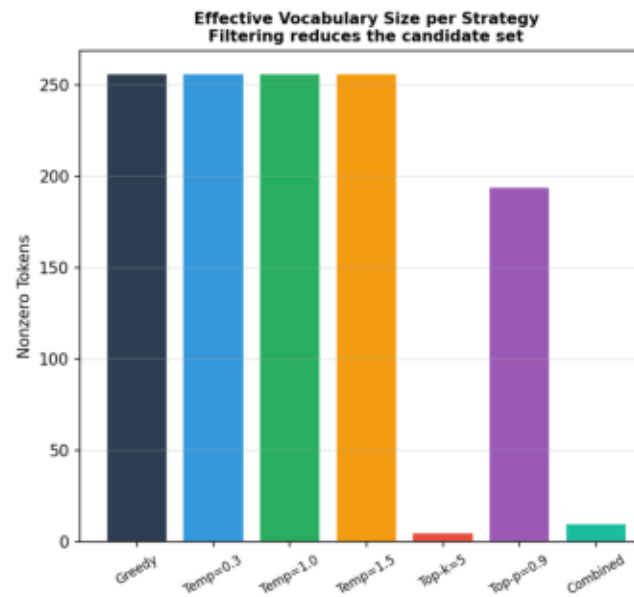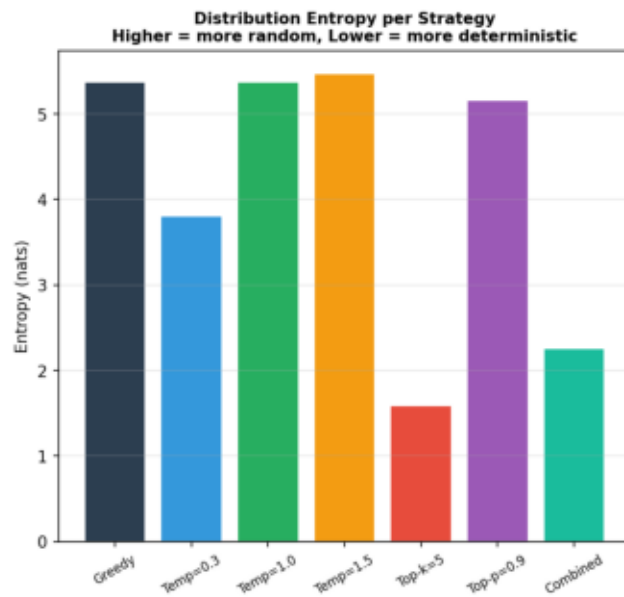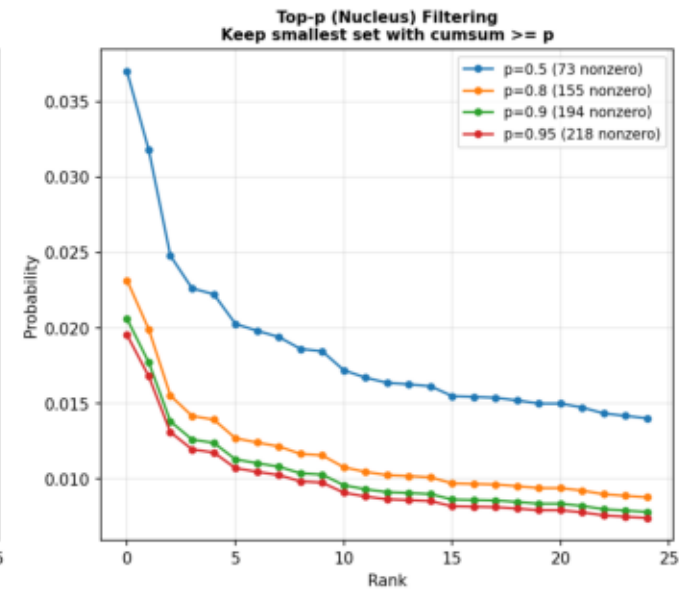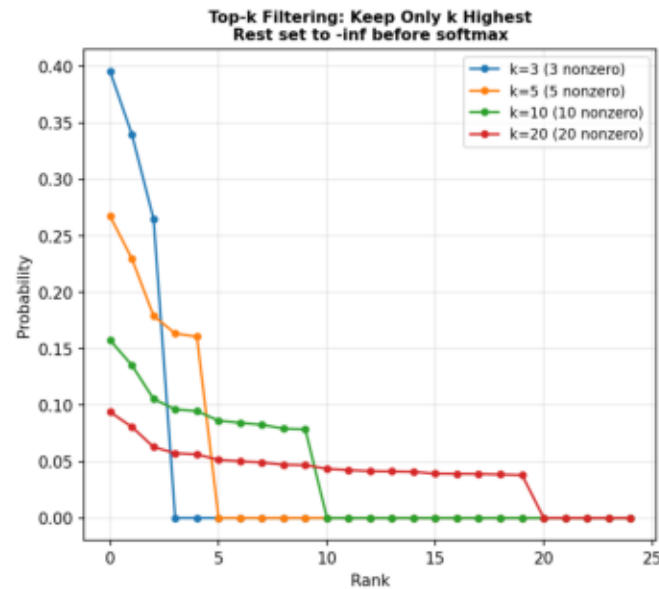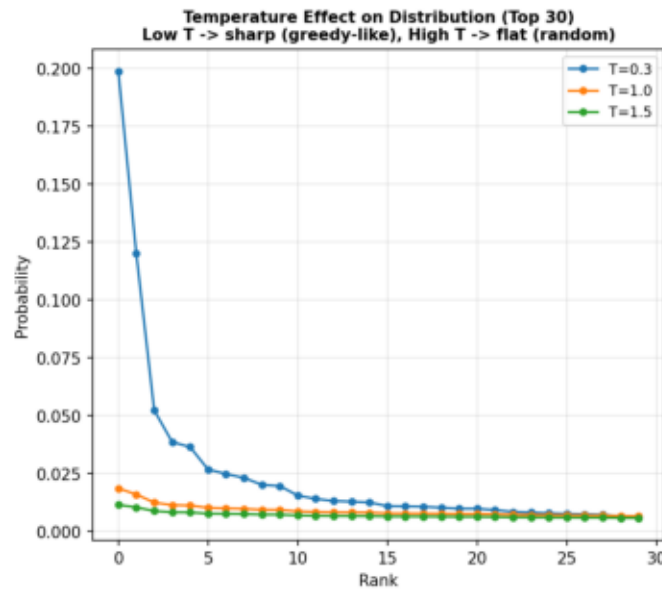
# Example 3: Sampling Strategy Comparison



Sampling Strategies: Temperature, Top-k, Top-p, and Combined

# Example 4: Autoregressive Generation Loop

## Autoregressive Generation: Step-by-Step Forward Pass with Growing Context



### Growing Context: Full Recompute at Each Step
Blue = existing context, Red = new token position

### Per-Step Cost: Naive vs KV Cache
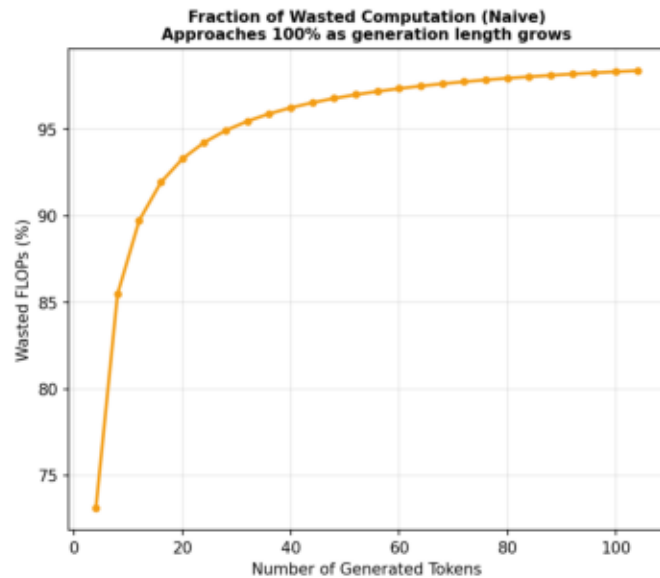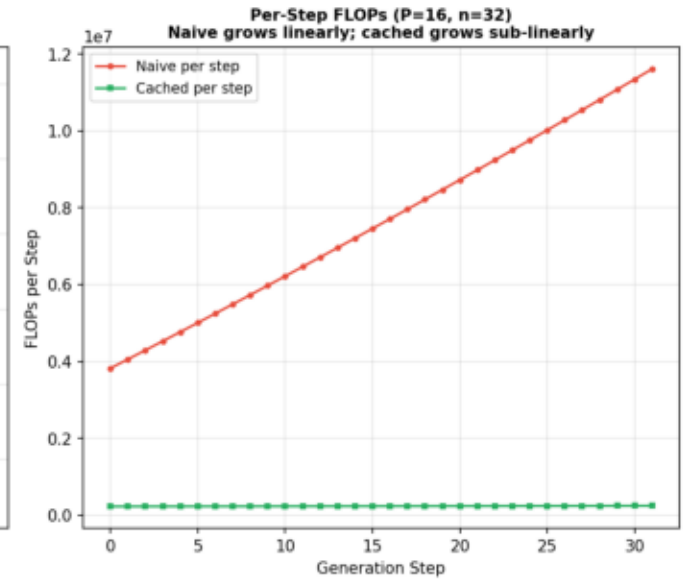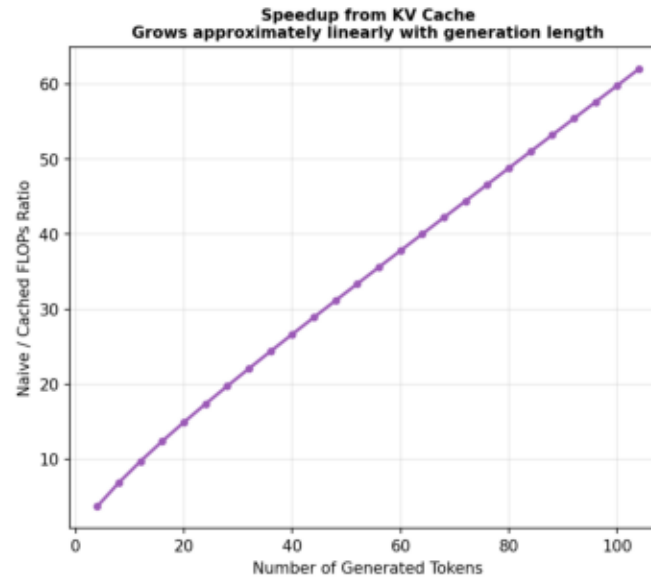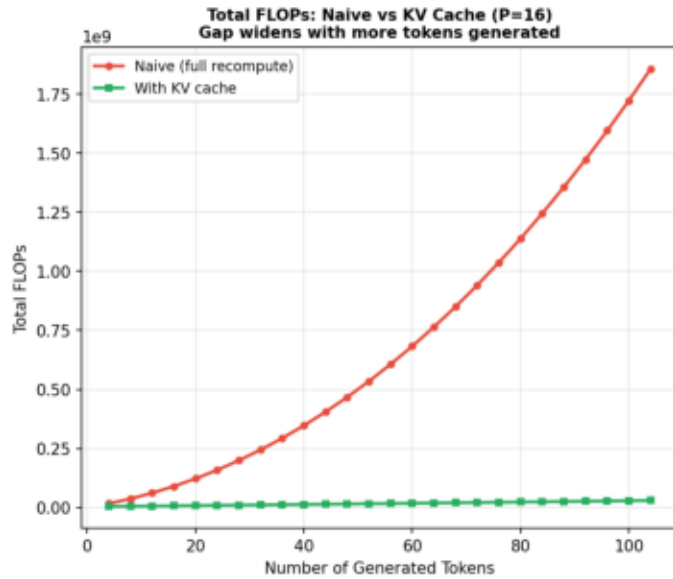Naive processes entire sequence; cache processes 1 token

- Naive (full recompute)
- With KV cache (1 token)

### Cumulative Cost: Quadratic vs Linear
After 10 steps: naive=85, cached=14

- Naive (O(n^2))
- Cached (O(n))
- Wasted computation

### Top-5 Token Probabilities per Step
(Random weights -> near-uniform distribution)

Rank 1, Rank 2, Rank 3, Rank 4, Rank 5

### Distribution Entropy per Step
Near-max entropy with random weights (expected)

Max entropy (uniform) = ln(256) = 5.55

```
AUTOREGRESSIVE GENERATION
=========================

Prompt: [3, 7, 11, 15]
Generated: [103, 103, 182, 117, 117, 117, 117, 117, 187, 252]

Algorithm (naive, no KV cache):
    for step in range(max_new):
        logits = model.forward(tokens)  # FULL
        next = sample(logits[:, -1, :]) # last
        tokens = concat(tokens, next)

Cost analysis:
    P=4, n=10
    Step i processes (P+i) tokens
    Total = n*P + n(n-1)/2
          = 85 token-steps
    With KV cache: P + n = 14
    Ratio: 6.1x

For large n, naive is O(n^2),
cached is O(n). This quadratic
waste motivates the KV cache.
```

# Example 5: Computational Cost (Naive vs KV Cache)



Computational Cost: Naive Full Recompute vs Theoretical KV Cache

# Example 6: Model Parameter Analysis



Model Parameter Analysis: Llama 2 7B, Llama 3 8B, Mistral 7B

**Total Parameter Count (no weight tying)**
**Transformer blocks dominate (>86%)**

- Embedding
- Blocks (N x block)
- Output Proj

**Llama 2 7B**
Total = 6.74B (untied weights)

**Llama 3 8B**
Total = 8.03B (untied weights)

**Mistral 7B**
Total = 7.24B (untied weights)

**Hypothetical Weight Tying Savings**
**Would save d_model x V parameters**

- Actual (no tying)
- Hypothetical (with tying)

```
MODEL PARAMETER ANALYSIS
========================

Parameter formula (with tying):
  P = V*d + N*P_block + d

Per-block (Llama 2 7B):
  Attn:  67.1M (33.2%)
  FFN:   135.3M (66.8%)
  Norms: 8192 (~0%)

Weight tying (hypothetical):
  W out = E.T (shared memory)
  Would save V*d parameters
  Llama 3: would save 525M (128K vocab)

Key observations:
  - Blocks are >86% of total params
    (96% for 32K vocab, 87% for 128K)
  - FFN is ~67% of each block (MHA)
    (~81% with GQA due to fewer K/V)
  - Larger vocab (Llama 3) increases
    embedding + output proj cost
  - These models do NOT use weight
    tying (separate output projection)
```