# Multi-Head Attention

## Comprehensive Demo and Analysis

Parallel attention heads with fused weight matrices,
reshape/transpose operations, and output projection.
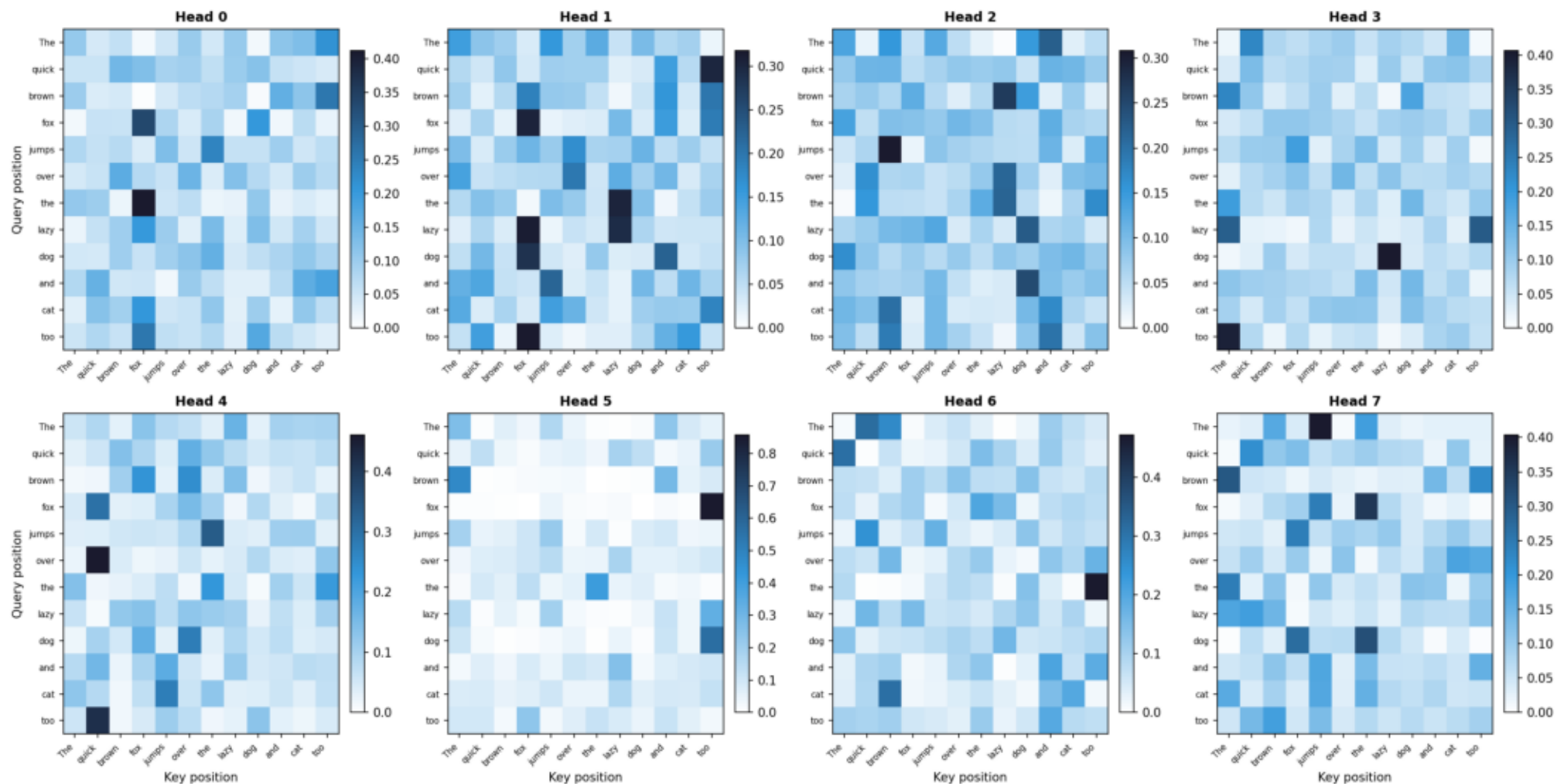
Random seed: 42
Number of visualizations: 7

# Summary of Findings

1. Attention Patterns: Each head learns distinct attention distributions, visible as different heatmap structures across 8 heads.

2. Multi-Head vs Single-Head: Multiple heads provide richer representational diversity. Standard deviation across heads reveals complementary patterns.

3. Head Diversity: Pairwise cosine similarity between head attention matrices confirms heads attend to different positions/relationships.

4. Causal Masking: Lower-triangular structure enforced correctly. Future positions receive zero attention weight. Row sums remain 1.0.

5. Memory Scaling: Attention matrix memory grows $O(L^2)$. At long sequences (L>1024), the attention matrix dominates total intermediate memory.

6. FLOPs Breakdown: Projection GEMMs dominate at short sequences. Attention core ($QK^T + AV$) overtakes at longer sequences.

7. Single-Head Equivalence: MHA(h=1) matches SelfAttention: PASS Both forward and backward pass agree to machine precision.
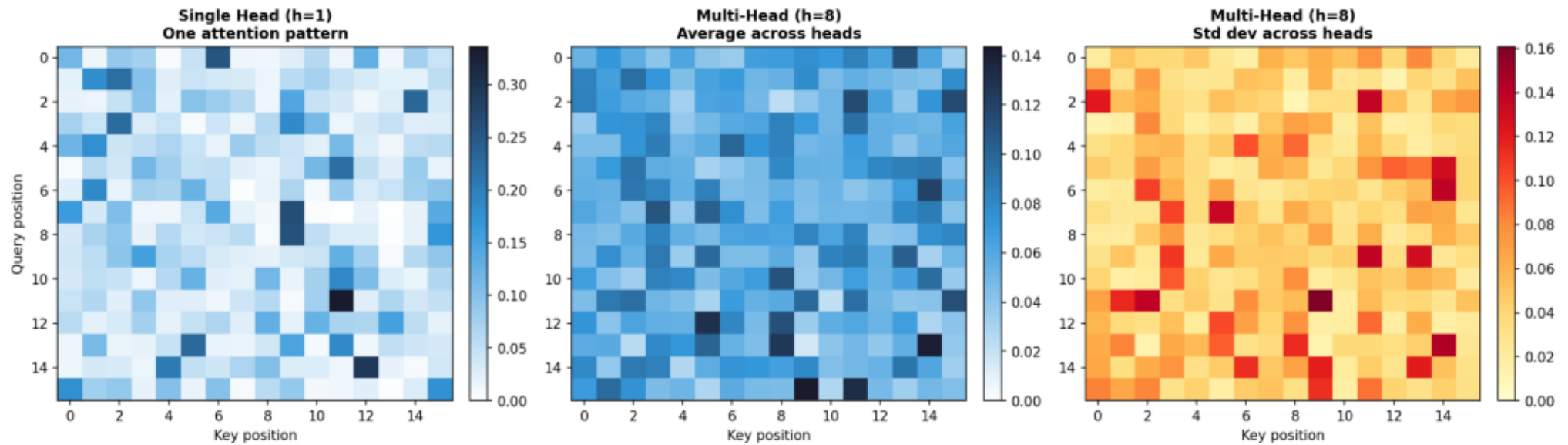
# Example 1: Per-Head Attention Heatmaps



Multi-Head Attention Patterns (8 Heads, d_model=64)

# Example 2: Multi-Head vs Single-Head Comparison

**Single-Head vs Multi-Head: Attention Diversity**



Single Head (h=1)
One attention pattern

Multi-Head (h=8)
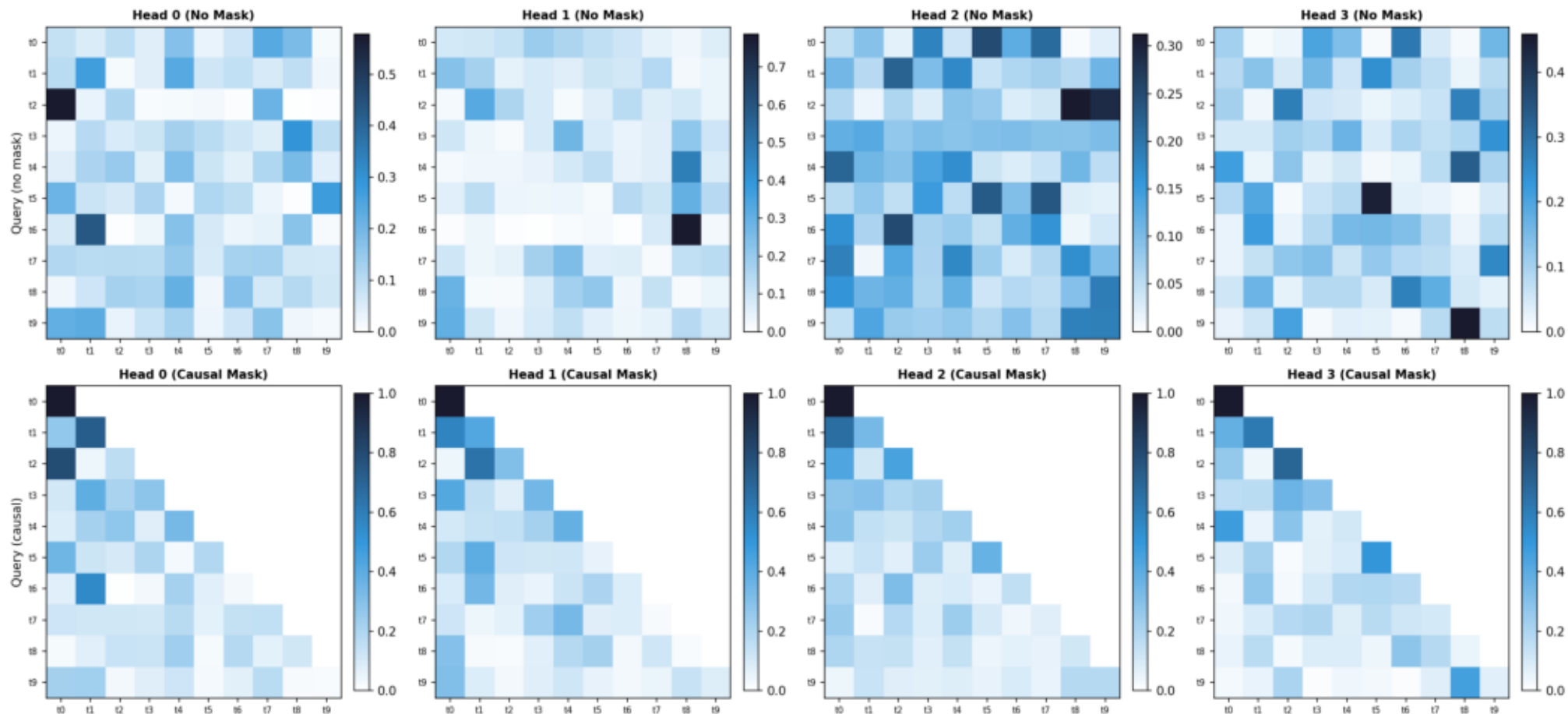Average across heads

Multi-Head (h=8)
Std dev across heads

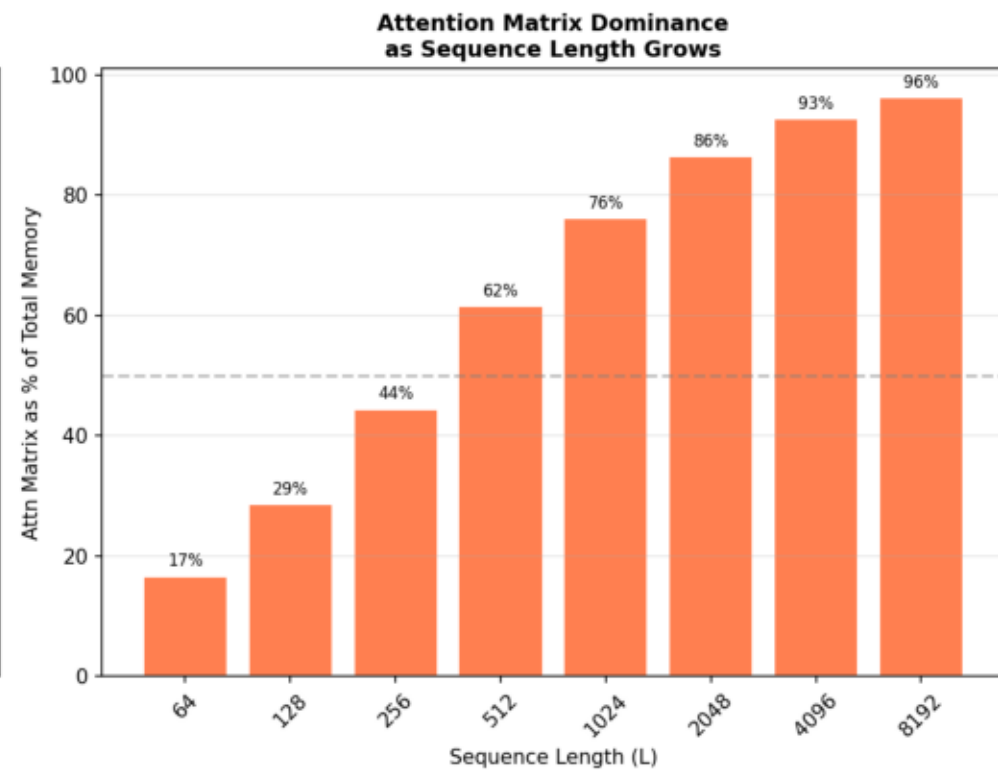**Example 3: Head Diversity Analysis (Cosine Similarity)**

# Example 4: Causal Masking Visualization



Effect of Causal Masking on Attention Patterns

# Example 5: Memory Scaling Analysis O(L^2)



**Intermediate Memory vs Sequence Length**
**(B=1, d=768, h=12)**

- Total FP32
- Total FP16
- Attn matrix only (FP32)

Memory (MB) vs Sequence Length (L): $2^6$, $2^7$, $2^8$, $2^9$, $2^{10}$, $2^{11}$, $2^{12}$, $2^{13}$

**Attention Matrix Dominance**
**as Sequence Length Grows**

Attn Matrix as % of Total Memory vs Sequence Length (L)

| L | % |
|------|-----|
| 64 | 17% |
| 128 | 29% |
| 256 | 44% |
| 512 | 62% |
| 1024 | 76% |
| 2048 | 86% |
| 4096 | 93% |
| 8192 | 96% |

# Example 6: FLOPs Breakdown by Component



**FLOPs Breakdown by Component**
**(B=1, d=768, h=12)**

- Projection GEMMs
- Attention Core (QK + AV)
- Softmax

GFLOPs vs Sequence Length (L): 32, 64, 128, 256, 512, 1024, 2048, 4096

**FLOPs Composition Shift**
**Projections vs Attention**

% of Total FLOPs vs Sequence Length (L): $2^5$, $2^6$, $2^7$, $2^8$, $2^9$, $2^{10}$, $2^{11}$, $2^{12}$

Analytical L=1536

- Projection GEMMs
- Attention Core
- Softmax

# Example 7: Single-Head Equivalence Verification



MHA (h=1) vs SelfAttention: Numerical Equivalence
(d_model=32, B=2, L=8)

| | Max Absolute Difference |
|---|---|
| Forward (no mask) | 0.00e+00 |
| Forward (causal mask) | 0.00e+00 |
| Backward dX | 0.00e+00 |
| Backward dW_Q | 0.00e+00 |
| Backward dW_K | 0.00e+00 |
| Backward dW_V | 0.00e+00 |
| Backward dW_O | 0.00e+00 |

--- Machine epsilon (~1e-12)