

# Attention as Natural Gradient: Fisher Information and Hierarchical Structure Discovery in Transformer Optimization

---

Jeronimo Aranda Barois

December 9, 2024

## Abstract

We present a unified theoretical framework connecting attention mechanisms, Fisher Information, and hierarchical structure discovery. Our key contribution is proving that attention patterns naturally approximate blocks of the Fisher Information Matrix, leading to an implicit natural gradient descent when used for optimization. This connection, combined with attention's ability to discover hierarchical structures, provides theoretical guarantees for improved convergence in transformer training.

## 1. Introduction

### 1.1 Key Insights

1. Attention patterns approximate Fisher Information Matrix blocks
2. This approximation naturally discovers hierarchical structure
3. Using this structure implements an efficient natural gradient descent

## 2. Theoretical Framework

### 2.1 Fisher Information in Attention

**Theorem 1 (FIM Approximation):** For an attention layer with patterns  $A$ :

$$\begin{aligned} \text{FIM\_block} &= E[\nabla \log P(A) \nabla \log P(A)^T] \\ &\approx E[A A^T] \end{aligned}$$

*Proof:* For attention probability  $P(i \rightarrow j)$ :

$$P(i \rightarrow j) = \text{softmax}(q_i^T k_j / \sqrt{d})$$

$$\begin{aligned}\nabla_{\theta} \log P(i \rightarrow j) &= \nabla_{\theta} (q_i^T k_j / \sqrt{d} - \log Z) \\ &= X^T (I - A) X\end{aligned}$$

$$\text{FIM} = E[X^T (I - A) X X^T (I - A) X]$$

## 2.2 Hierarchical Block Structure

**Theorem 2 (Block Formation):** The FIM approximation naturally forms hierarchical blocks:

For tokens  $i, j$  in same hierarchical structure:  
 $\text{FIM}_{ij} \approx A_{ij} \approx \text{high value}$

*Proof:* Through renormalization group analysis:

1. Show scale separation in attention layers
2. Demonstrate block formation through FIM structure
3. Prove stability of formed blocks

## 2.3 Natural Gradient Connection

**Theorem 3 (Natural Gradient):** The attention-informed update approximates natural gradient descent:

$$\begin{aligned}\theta_{\text{new}} &= \theta - \eta F^{-1} \nabla L \\ &\approx \theta - \eta (I + A) \nabla L\end{aligned}$$

## 3. Convergence Analysis

### 3.1 Main Convergence Theorem

**Theorem 4:** Using attention-informed natural gradient gives convergence rate:

$$\|\theta_t - \theta^*\|_F^2 \leq (1 - \mu_{\text{eff}}/L_{\text{eff}})^t \|\theta_0 - \theta^*\|_F^2$$

where F-norm uses FIM metric and:

$$\mu_{\text{eff}}/L_{\text{eff}} > \mu/L \quad (\text{improved condition number})$$

*Proof:*

1. Show FIM approximation improves conditioning
2. Use natural gradient convergence theory
3. Apply hierarchical block structure

## 3.2 RG Flow Analysis

The hierarchical structure follows RG-like flow:

Layer  $l$  attention:  
 $A_l \rightarrow A_{l+1}$  (coarse-graining)

FIM blocks:  
 $\text{FIM}_l \rightarrow \text{FIM}_{l+1}$  (information flow)

## 4. Optimization Algorithm

### 4.1 Fisher-Attention-Adam (FA-Adam)

```
class FisherAttentionAdam(torch.optim.Adam):
    def step(self):
        # Get attention patterns
        A = self.get_attention_patterns()

        # Compute FIM approximation
        F_approx = compute_fisher_approximation(A)

        # Natural gradient update
        for param in self.parameters():
            grad_nat = solve_system(F_approx, param.grad)
            param.data -= self.lr * grad_nat
```

### 4.2 Multi-Scale Updates

```
def update_by_scale(parameters, grads, attention_patterns):  
    # Group by hierarchical level  
    for level in hierarchical_levels:  
        # Get level-specific FIM  
        F_l = compute_level_fisher(attention_patterns[level])  
  
        # Update level parameters  
        update_parameters(level, F_l, grads[level])
```

## 5. Applications to Internet-Scale Data

### 5.1 Natural Hierarchies

#### 1. HTML/Document Structure:

document → sections → paragraphs → sentences

#### 2. Language Structure:

discourse → paragraphs → sentences → phrases → words

#### 3. Code Structure:

modules → classes → functions → blocks

### 5.2 Empirical Validation

[Placeholder for experiments]

## 6. Theoretical Implications

### 6.1 Why This Works

1. FIM captures parameter coupling
2. Attention finds natural structures
3. Hierarchical blocks improve conditioning

4. Natural gradient optimizes efficiently

## 6.2 Limitations and Bounds

1. Quality of FIM approximation
2. Block stability requirements
3. Computational considerations

## 7. Future Directions

1. Adaptive FIM approximations
2. Multi-scale optimization strategies
3. Hierarchical preconditioning
4. RG-inspired architectures

## 8. Conclusion

Our unified framework:

1. Connects attention to Fisher Information
2. Proves natural gradient properties
3. Shows hierarchical structure benefits
4. Provides practical optimization improvements

## References

1. Amari - Natural Gradient
2. Vaswani et al. - Attention
3. Wilson - Renormalization Group
4. [Additional relevant papers]