

Transformers: Understanding the Pipeline

Input → Computation → Output → WHY

Week 5: Transformers

Complete Example: Predicting the Next Word

INPUT: "The cat sat on the ___"

COMPUTATION (6 Steps):

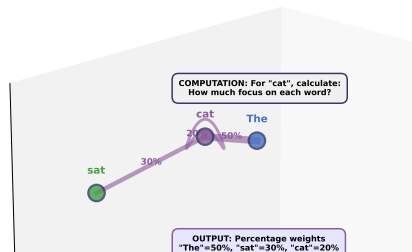
- ① Words \rightarrow Numbers [0.2, 0.5, ...]
- ② Add position: word 1, 2, 3, 4, 5
- ③ **Attention:** Each word looks at context
 - "sat" needs "cat" (WHAT sat?)
 - "on" needs "sat" (sitting ON something)
- ④ **Multi-Head Attention:**
 - Head 1: Grammar ("on the" \rightarrow needs noun)
 - Head 2: Meaning (cat + sat \rightarrow furniture)
 - Head 3: Position (final word prediction)
 - Head 4: Relations (cat sits ON things)
- ⑤ Combine all 4 perspectives
- ⑥ Predict next word

OUTPUT: Top predictions

- "mat": 85%
- "floor": 10%
- "table": 3%
- "rug": 2%

Result: "The cat sat on the **mat**"

Step 3: Calculate Attention (Who Looks at Who)



The Simple Goal

INPUT:

- Text: "The cat sat on the mat"
- 7 words (English)

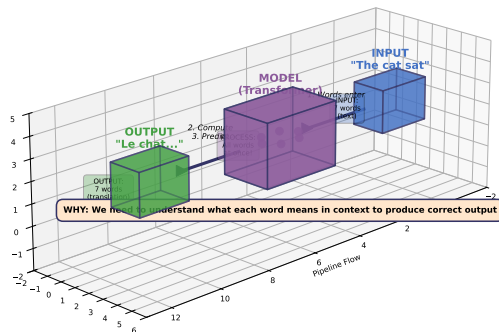
OUTPUT:

- Text: "Le chat était assis sur le tapis"
- 7 words (French)

THE TASK:

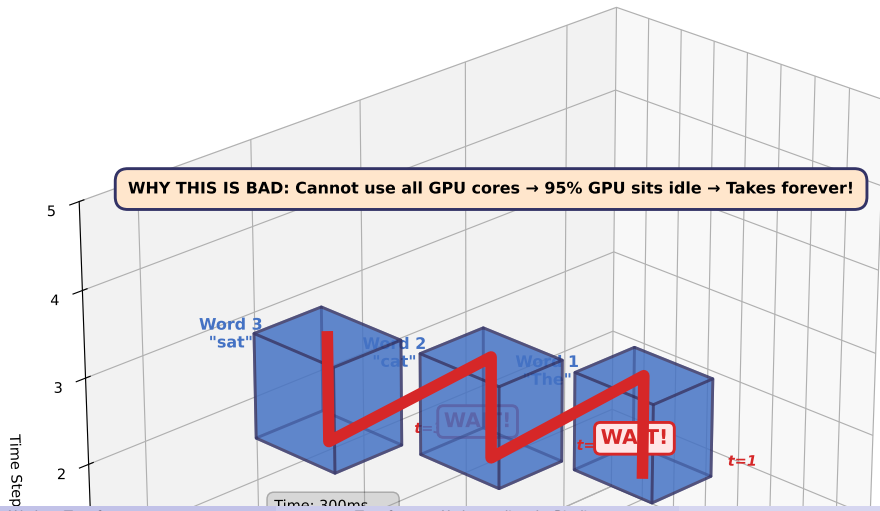
- Translate
- Predict next word
- Answer questions

The Transformer Pipeline: Input → Process → Output



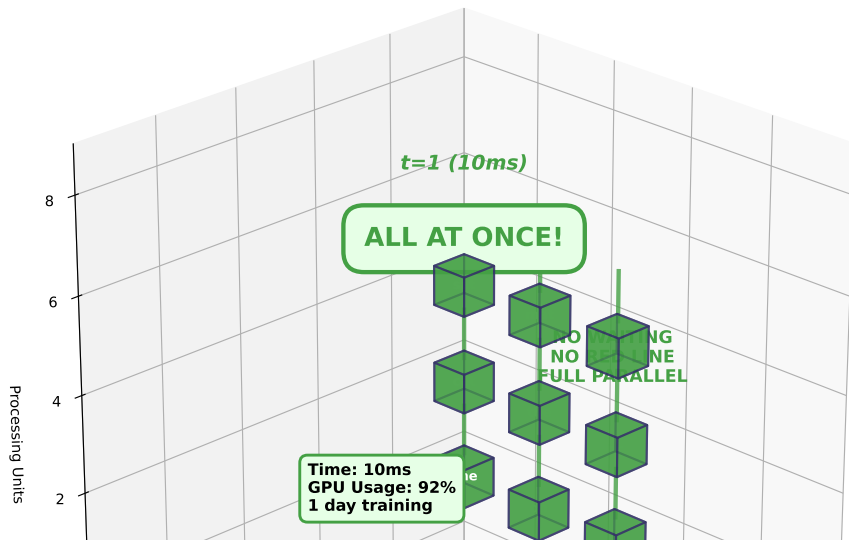
The Old Way: RNN (Sequential Processing)

RNN: Sequential Processing = RED LINE Bottleneck

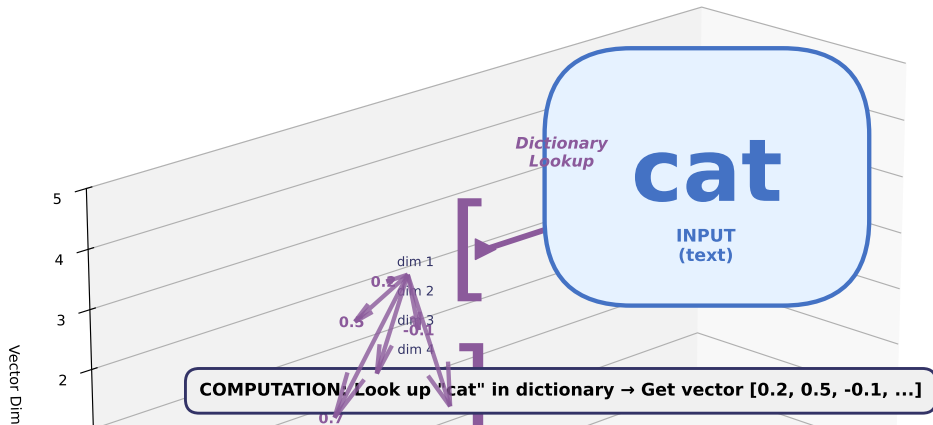


The New Way: Transformer (Parallel Processing)

Transformer: Parallel Processing = NO RED LINE!

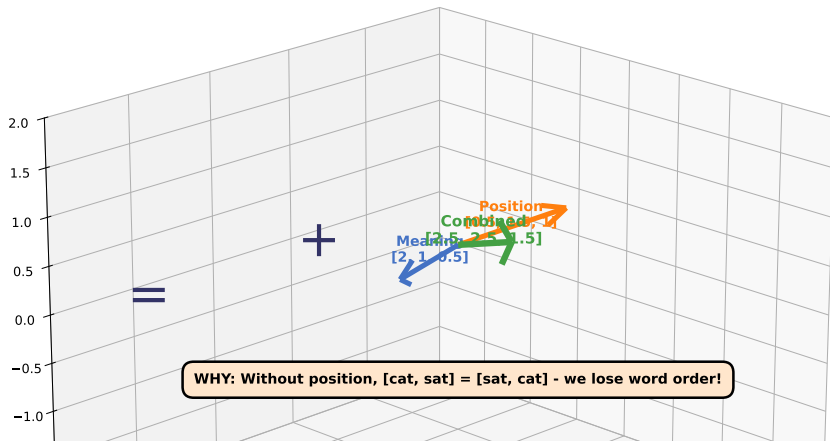


Step 1: Turn Words into Numbers



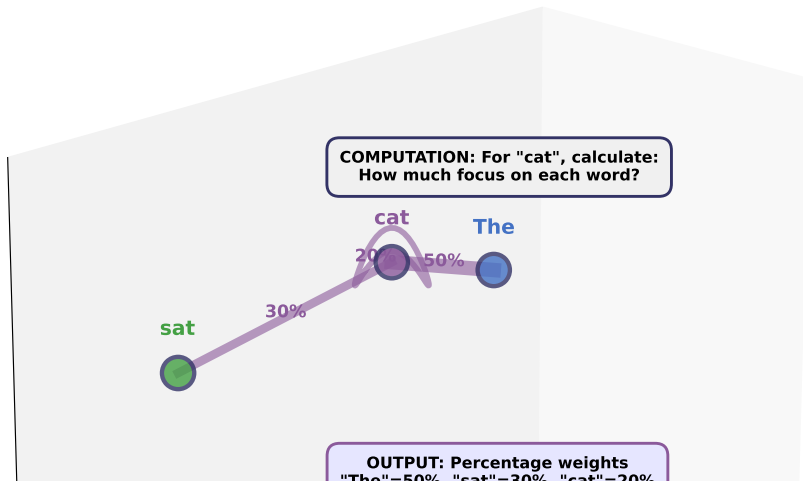
Step 2: Add Position Information

Step 2: Add Position Information



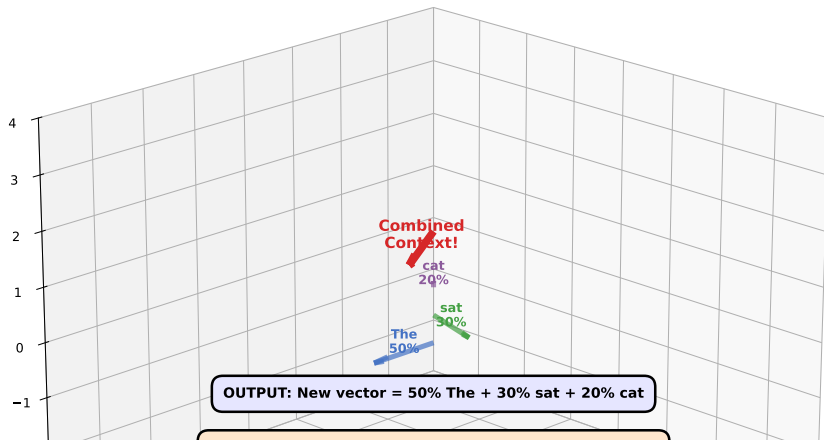
Step 3: Calculate Attention (Who Looks at Who)

Step 3: Calculate Attention (Who Looks at Who)



Step 4: Combine Information (Weighted Average)

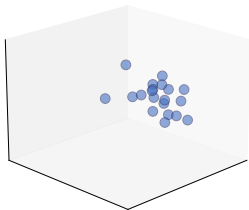
Step 4: Combine Information (Weighted Average)



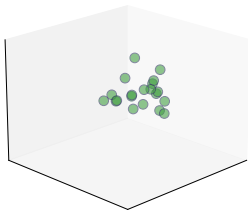
Step 5: Multiple Perspectives (Multi-Head Attention)

Step 5: Multiple Perspectives (8 Heads in Parallel)

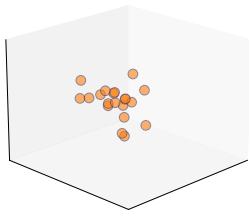
Head 1
Grammar



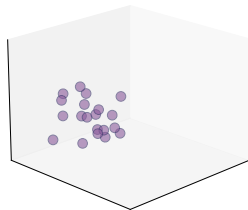
Head 2
Meaning



Head 3
Position



Head 4
Global



Step 6: Final Prediction

INPUT: Context-enriched vectors

- Each word knows about:
 - Its meaning
 - Its position
 - Related words (8 perspectives)

COMPUTATION:

- Feed through prediction layer
- Calculate probabilities for each possible next word

OUTPUT:

- Next word probabilities:
 - "Le": 85%
 - "The": 10%
 - Other: 5%
- Pick highest: "Le"

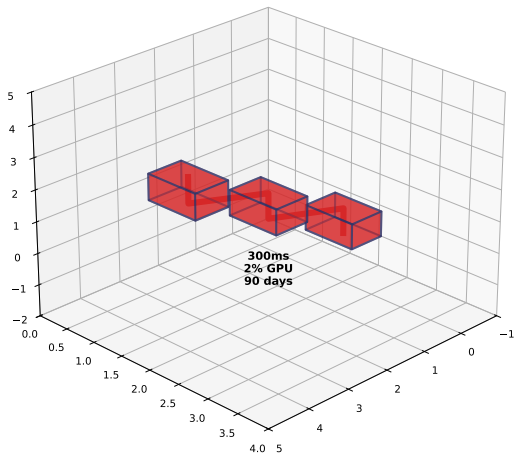
Result: Translation complete!

WHY: This is what we wanted all along - accurate prediction from context!

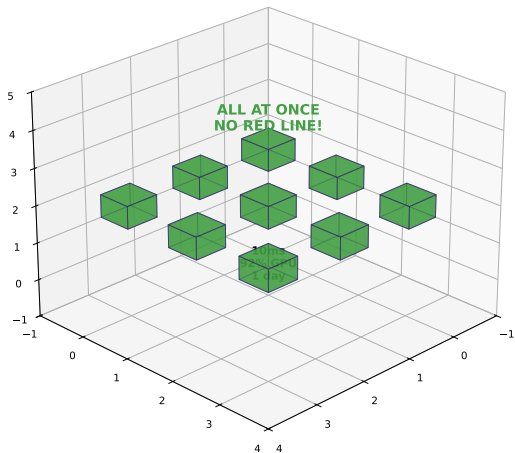
The Speed Secret: Parallel Processing

Speed Comparison: Sequential Staircase vs Parallel Cube

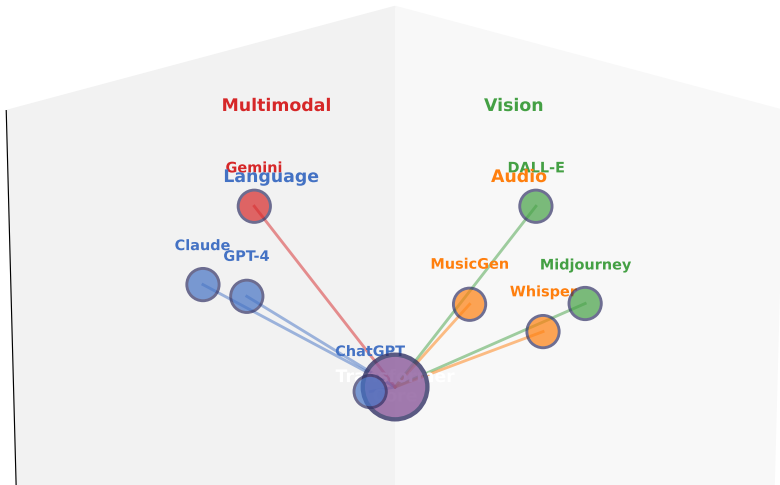
RNN: Sequential (RED LINE)



Transformer: Parallel (NO RED LINE)



2024 Landscape: Transformers Power Everything



The Tradeoff: What We Gave Up

Advantages (PRO):

- 100x faster training
- Parallel processing
- 92% GPU utilization
- Works on any data type
- Enabled modern AI

Disadvantages (CON):

- More memory (quadratic)
- Needs more data
- Limited sequence length
- More complex to tune

THE DECISION: Speed + quality $\hat{=}$ memory cost for modern AI

WHY ACCEPT TRADEOFF: Memory is cheap, time is expensive. Better to train fast even if uses more RAM.

Summary: The Pipeline Recap

The 6-Step Pipeline:

- 1 **Words → Numbers:** Dictionary lookup (embeddings)
- 2 **Add Position:** Vector addition (meaning + position)
- 3 **Calculate Attention:** Who looks at who? (percentage weights)
- 4 **Combine Information:** Weighted average (context-enriched)
- 5 **Multiple Perspectives:** 8 heads in parallel (grammar, meaning, position, ...)
- 6 **Predict Output:** Final layer (translation/next word)

KEY INSIGHT: All words processed in parallel!

- Result: 90 days → 1 day (90x speedup)
- Enabled: ChatGPT, GPT-4, DALL-E, Whisper, ...

Next Week: Pre-training & Fine-tuning - Now that training is fast, we can train HUGE models!

Transformers

Understanding the Pipeline

Input → Computation → Output → WHY

Questions?