# Transformers: Understanding the Pipeline
Input $\rightarrow$ Computation $\rightarrow$ Output $\rightarrow$ WHY

Week 5: Transformers

# The Simple Goal

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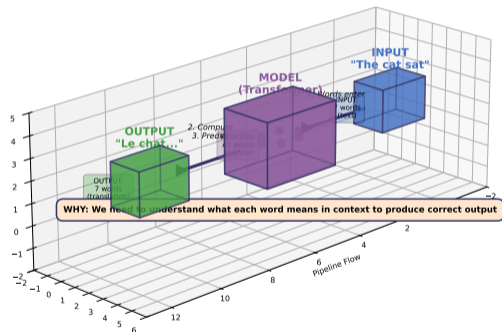- Text: "The cat sat on the mat"
- 7 words (English)

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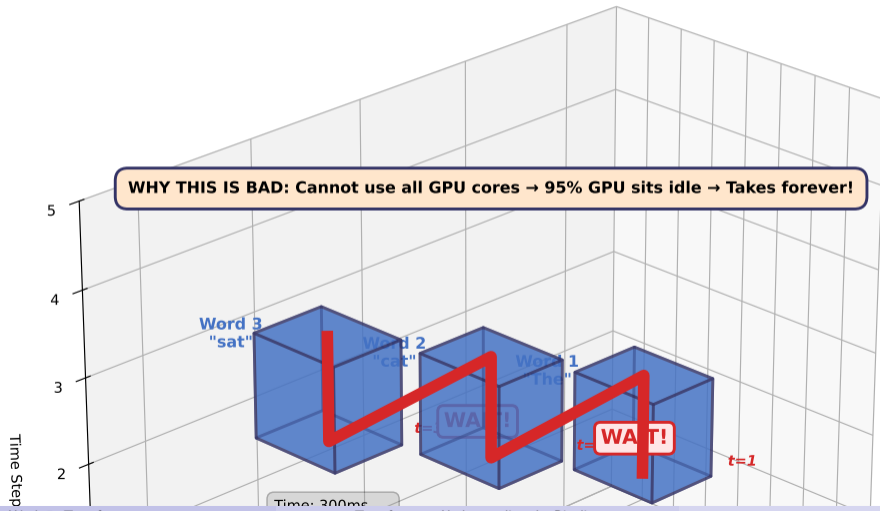- Text: "Le chat était assis sur le tapis"
- 7 words (French)

**THE TASK:**
- Translate
- Predict next word
- Answer questions



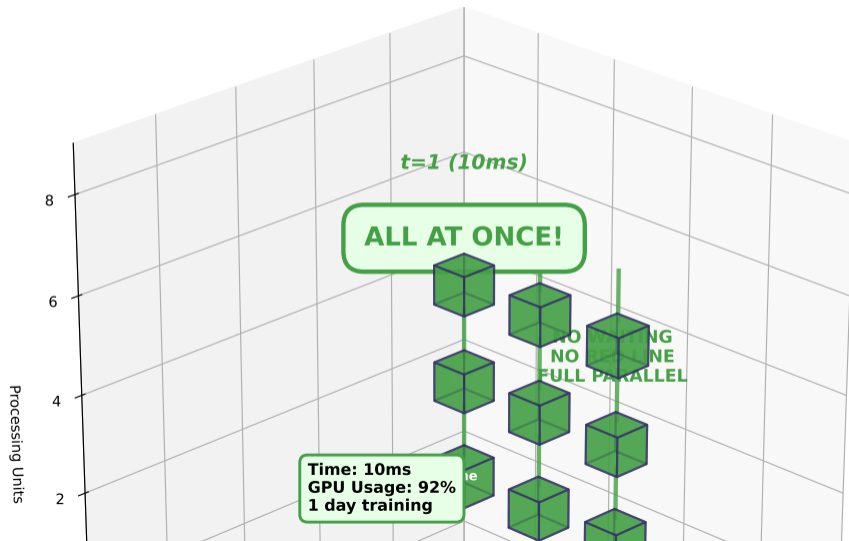The Transformer Pipeline: Input → Process → Output

**RNN: Sequential Processing = RED LINE Bottleneck**



WHY THIS IS BAD: Cannot use all GPU cores → 95% GPU sits idle → Takes forever!

**Transformer: Parallel Processing = NO RED LINE!**



*t=1 (10ms)*

**ALL AT ONCE!**

NO WAITING
NO RED LINE
FULL PARALLEL

Time: 10ms
GPU Usage: 92%
1 day training

Processing Units

**Step 1: Turn Words into Numbers**



*Dictionary Lookup*

**cat**

**INPUT
(text)**

dim 1
**0.2**
dim 2
**0.5**
dim 3
**-0.1**
dim 4
**0.7**

Vector Dim

COMPUTATION: Look up "cat" in dictionary → Get vector [0.2, 0.5, -0.1, ...]

## Step 2: Add Position Information



Position

Combined: [2.5, 1.5]

Meaning: 2.5, 1.5]

[2, 1, 5.5]

+

=

WHY: Without position, [cat, sat] = [sat, cat] - we lose word order!

## Step 3: Calculate Attention (Who Looks at Who)



COMPUTATION: For "cat", calculate:
How much focus on each word?

cat

The

20%   50%

sat

30%

OUTPUT: Percentage weights
"The"=50%, "cat"=30%, "cat"=20%

**Step 4: Combine Information (Weighted Average)**



OUTPUT: New vector = 50% The + 30% sat + 20% cat

**Step 5: Multiple Perspectives (8 Heads in Parallel)**
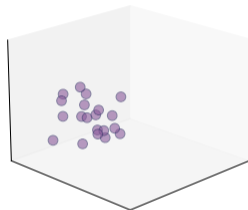


**Head 1**
**Grammar**

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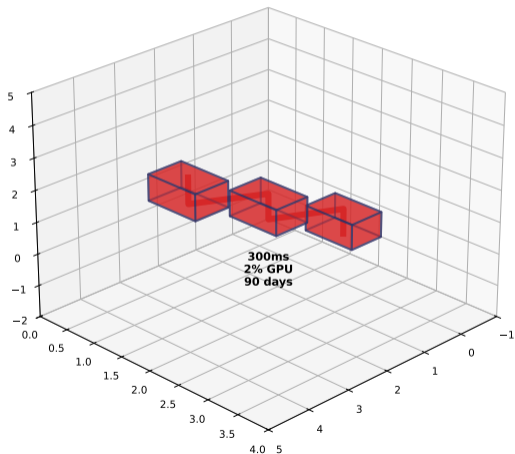- Next word probabilities:
  - "Le": 85%
  - "The": 10%
  - Other: 5%
- Pick highest: "Le"

**Result:** Translation complete!

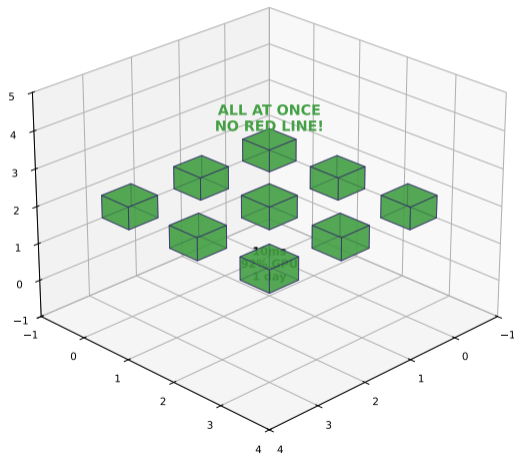**WHY:** This is what we wanted all along - accurate prediction from context!

Speed Comparison: Sequential Staircase vs Parallel Cube

RNN: Sequential (RED LINE)

Transformer: Parallel (NO RED LINE)

300ms
2% GPU
90 days

ALL AT ONCE
NO RED LINE!

# Real Numbers: The Proof

**Actual Experimental Results ("Attention Is All You Need", 2017):**

| Model | Training Time | GPU Usage | Quality (BLEU) |
|-------|---------------|-----------|----------------|
| RNN | 90 days | 2% | 24.5 |
| RNN + Attention | 45 days | 5% | 28.4 |
| **Transformer** | **1 day** | **92%** | **28.4** |

**KEY INSIGHT:**
- Same quality in 45x less time
- Better GPU utilization (2% $\rightarrow$ 92%)
- Enabled modern AI scale

**WHY BELIEVE THIS:** Published results, reproduced worldwide, powers all modern LLMs

**2024 Landscape: Transformers Power Everything**

# The Tradeoff: What We Gave Up

**Advantages (PRO):**

- 100x faster training
- No sequential bottleneck
- 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 ¿ 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 (NO RED LINE!):**

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 in parallel - NO RED LINE!**

- 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 $\rightarrow$ Computation $\rightarrow$ Output $\rightarrow$ WHY

Questions?