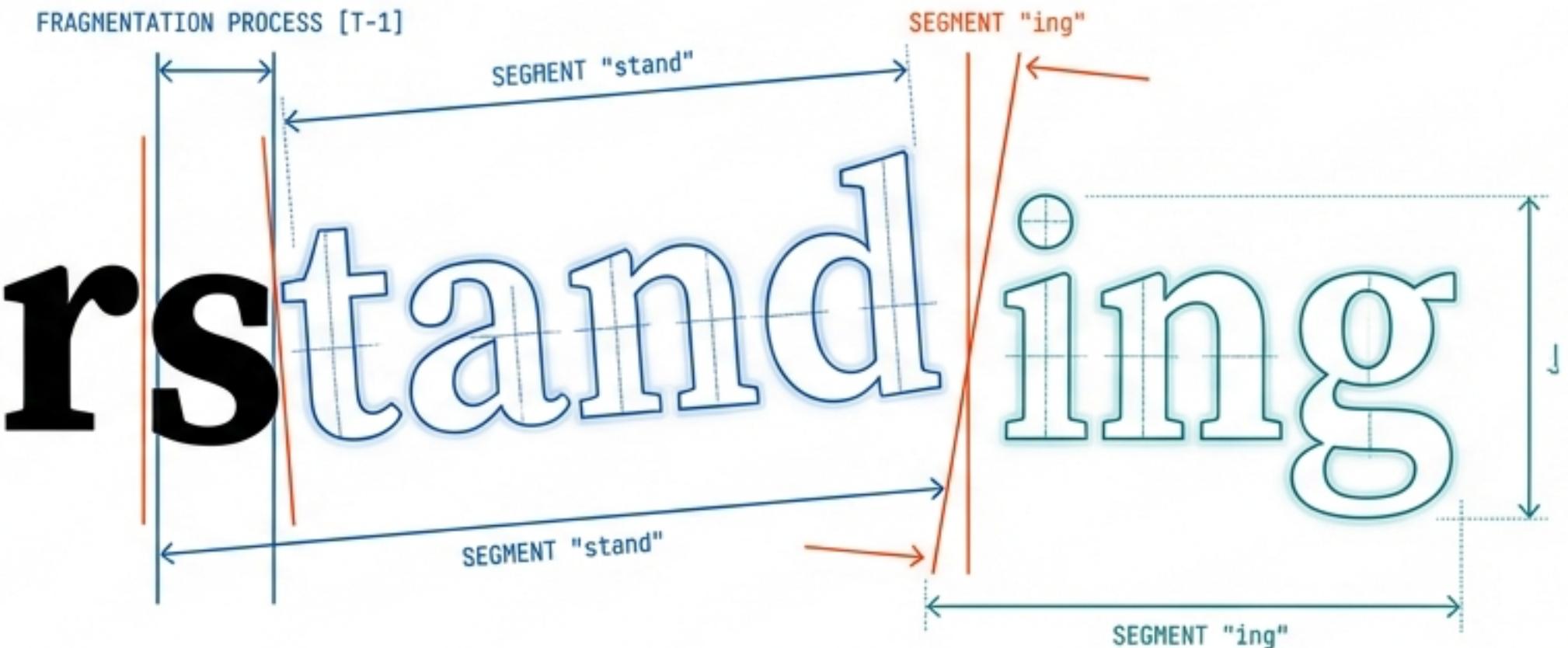


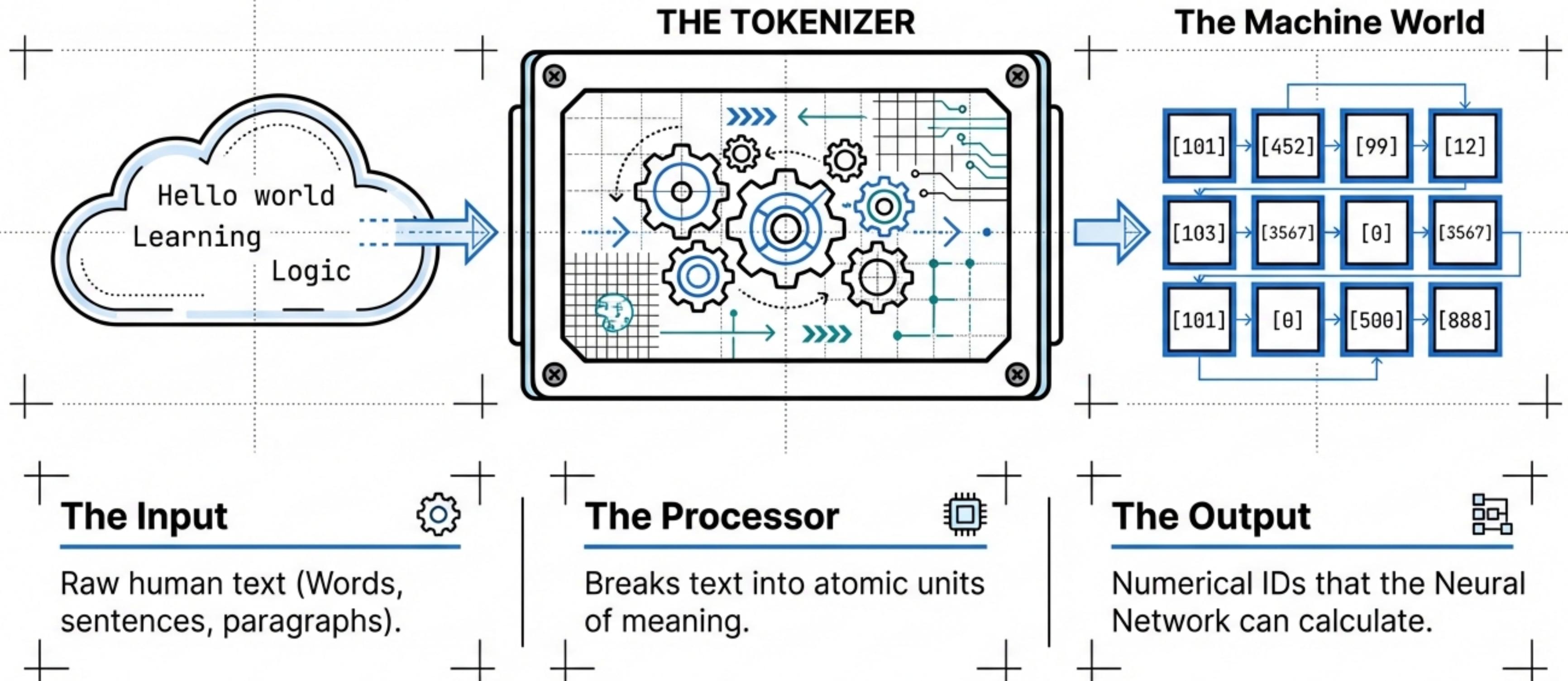
# Understand1ing



## The Mechanics of Subword Tokenization

How modern AI translates human language into machine logic.

# The Bridge Between Worlds



## Approach 1: The Trap of Full Word Tokenization

Treating every variation as a unique object.

playing player playful

Alert Orange
[playing]

Alert Orange
[player]

Alert Orange
[playful]

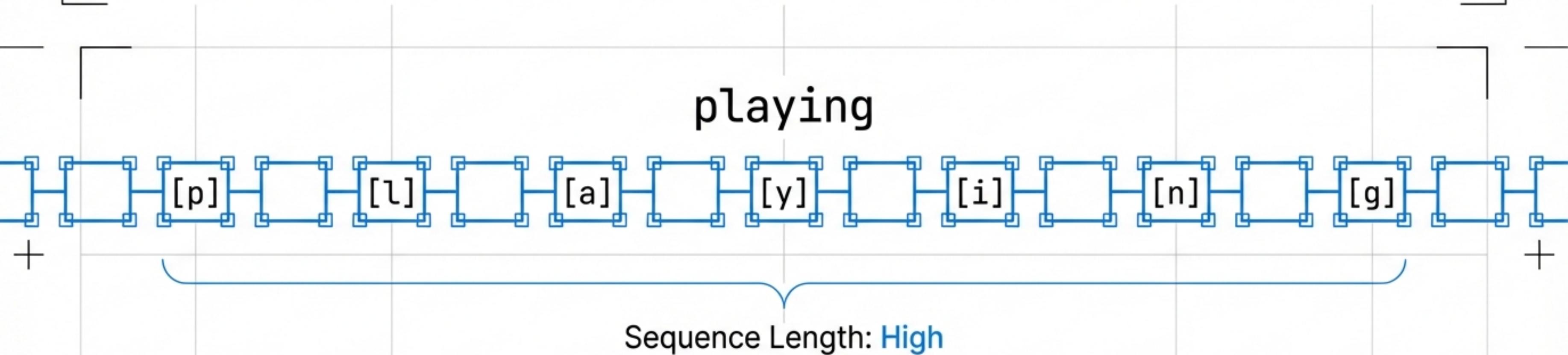
The Vocabulary Explosion Problem:

- 1. No shared meaning: The model does not know these words share the root 'play'.
- 2. Massive Size: English has >170,000 words. Storing all variations is inefficient.
- 3. Fragility: If the model sees 'played' but it's not in the list, it crashes (Unknown Token).



# Approach 2: The Inefficiency of Character Tokenization

Breaking structure down to dust.



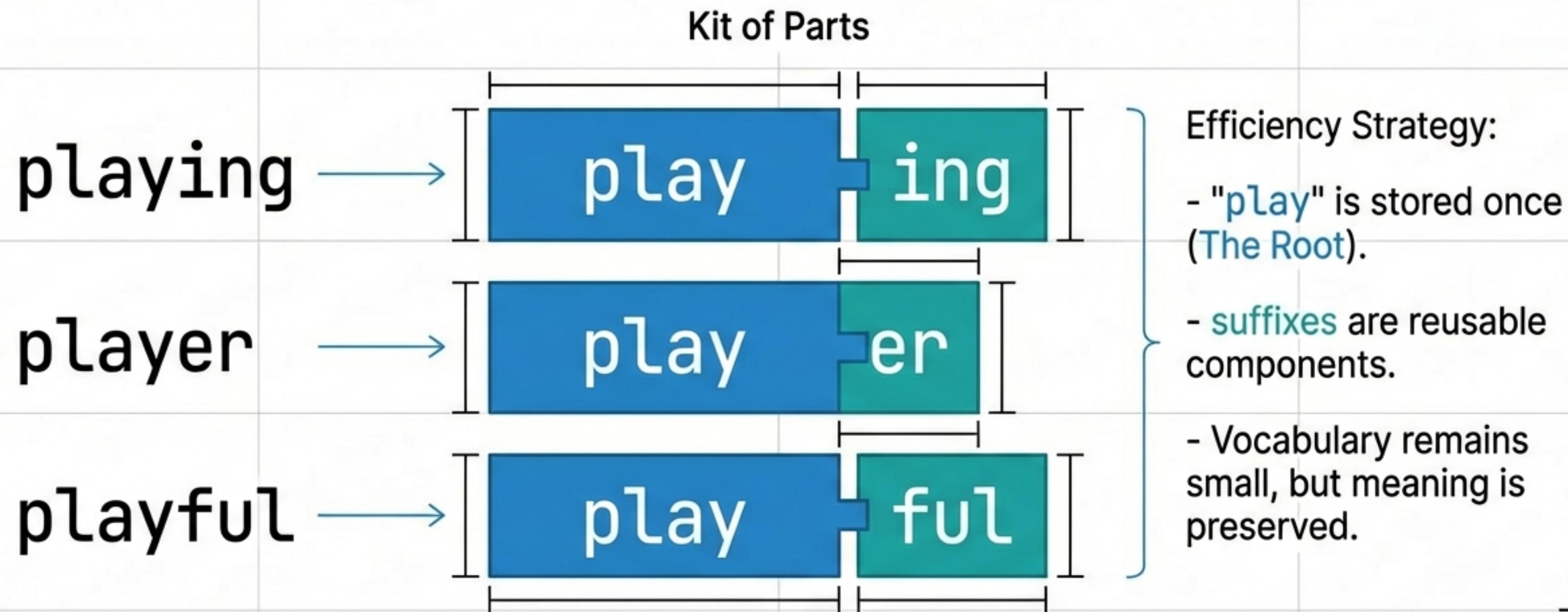
## The Sequence Overload Problem:

- 1. Loss of Meaning: The letter 'p' has no semantic meaning on its own.
- 2. Computational Cost: Processing thousands of characters for simple sentences is slow.
- 3. Training Difficulty: The model must learn how to spell before it can learn how to think.



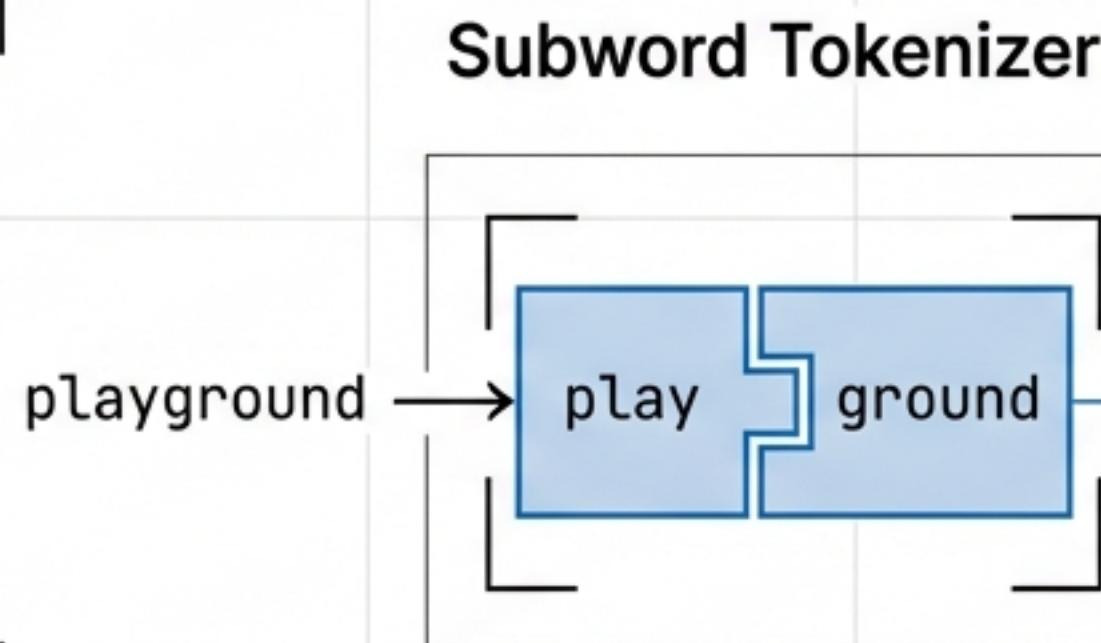
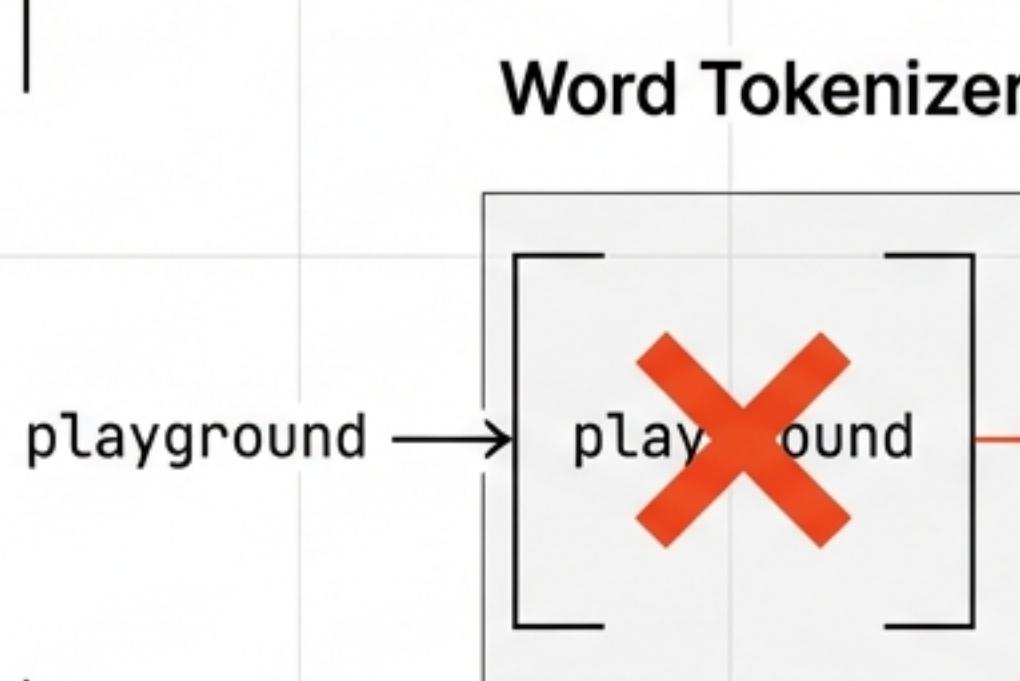
# The Solution: Subword Tokenization

The Goldilocks Zone: Semantic roots + reusable modifiers.



# The Power of Generalization

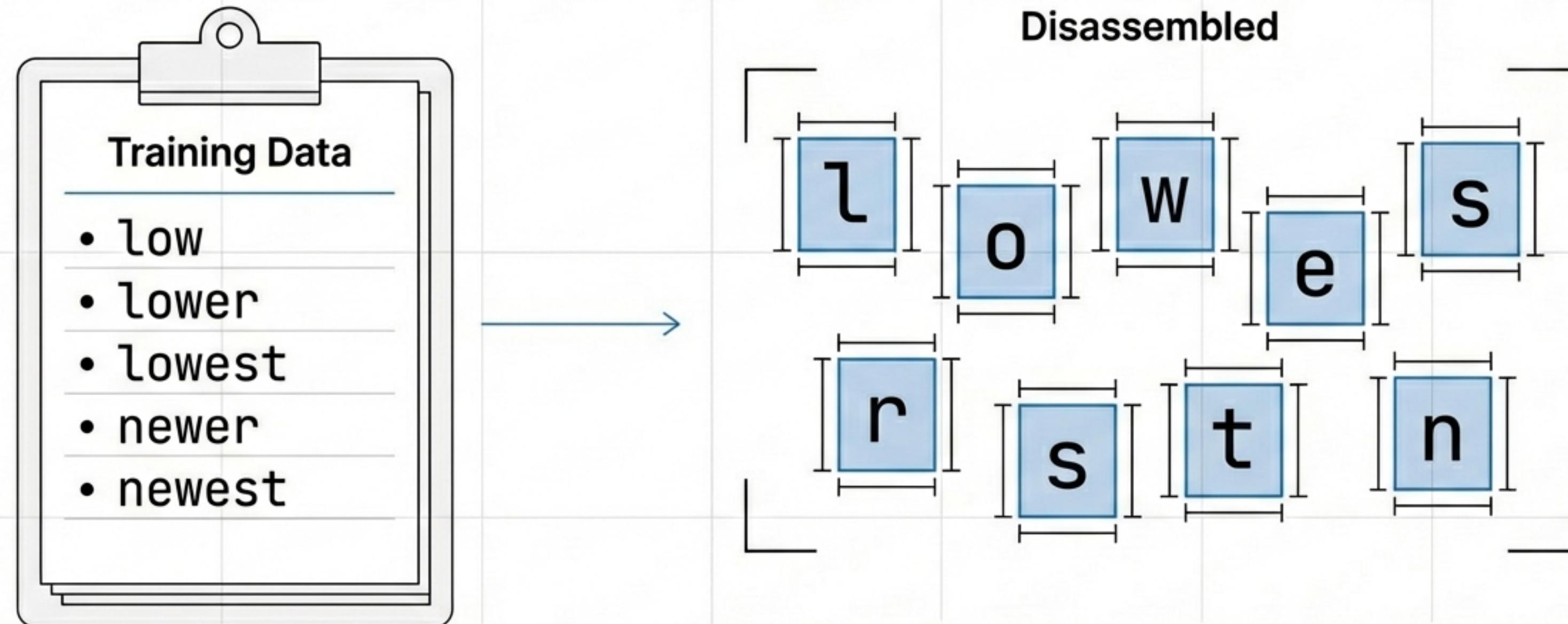
Understanding words the model has never seen before.



The model understands “playground” is a location for playing, simply by combining existing knowledge.

# Building the Tokenizer: The Training Phase

Step 1: Analyzing the Raw Corpus.



# The Merge: Learning Patterns

Algorithms (like BPE) iteratively merge frequent pairs.

Step 1 (Pair Found)

l + o



lo

Frequent Pair



Step 2 (Expansion)

lo + w



low

New Token  
Created

Step 3 (Parallel Learning)

e + s + t



est

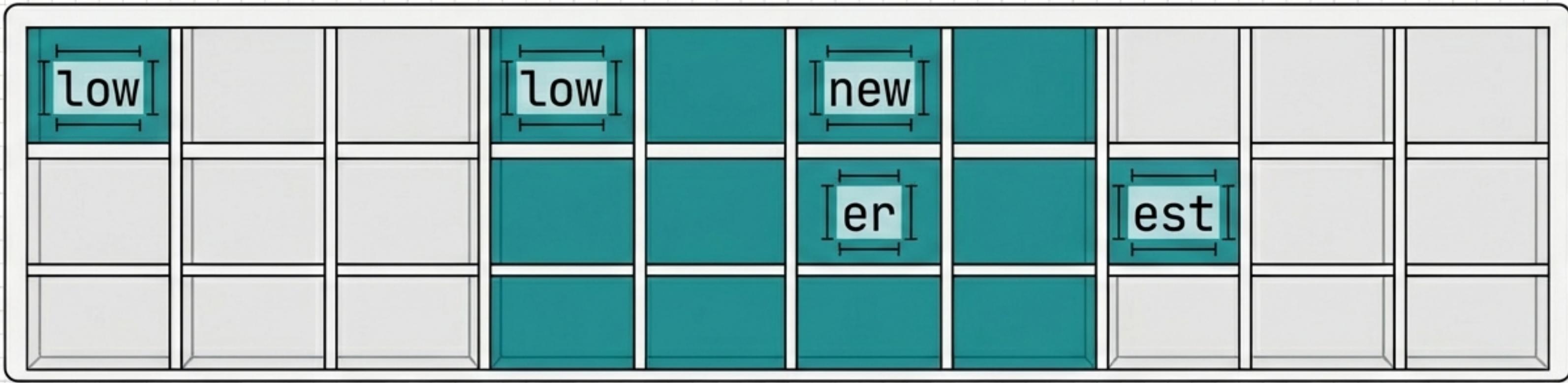
n + e + w



new

# The Optimized Vocabulary

The Final Toolset.



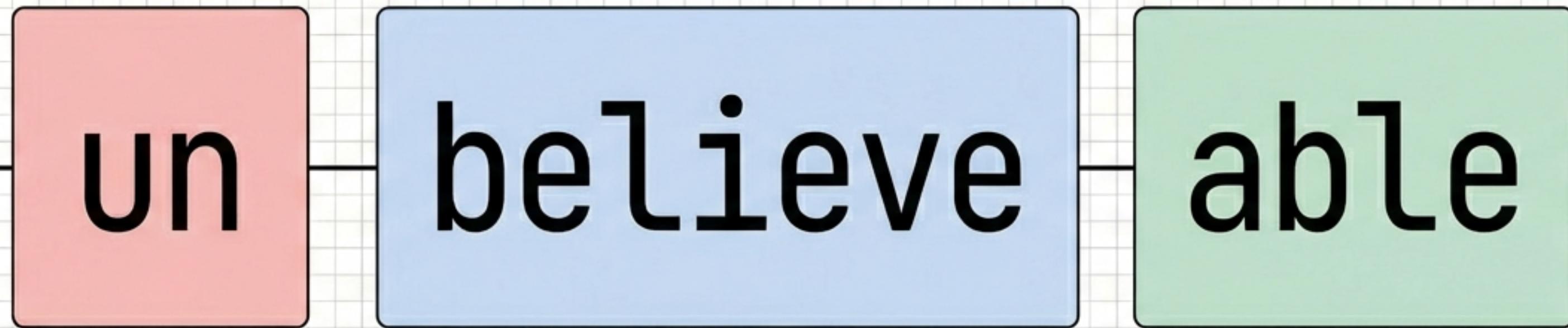
lower = low + er

lowest = low + est

newest = new + est

**100% Coverage with Minimal Storage.**

# Handling Complexity: The BERT Example



Prefix  
(Negation)

Root  
(Action/Concept)

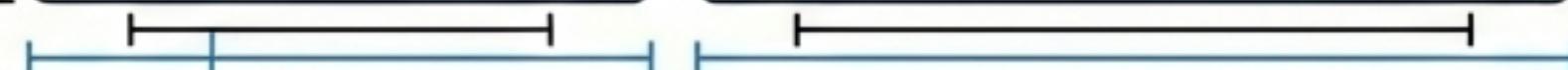
Suffix  
(Possibility)

Even if 'unbelievable' is rare in the data, the model understands:  
Not (un) + Believe + Able.

# Modern Tokenization: The GPT Approach

Statistical splitting over strict linguistics.

international → [ [inter] [national] ]



Clean Semantic Split

chatgpt → [chat] [g] [pt]

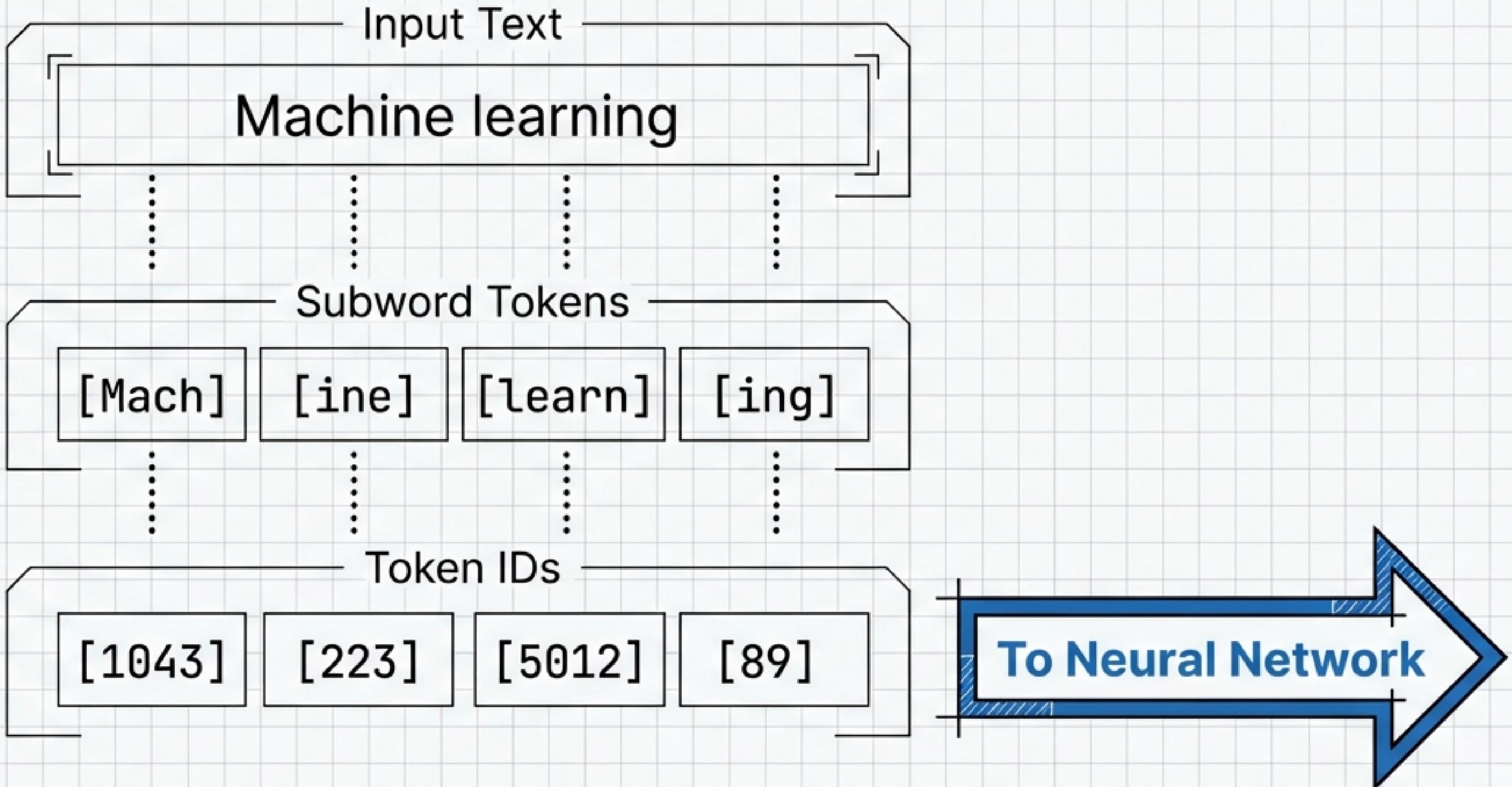
Frequency-based Split (Jargon)

\_ hello → [ \_hello ]

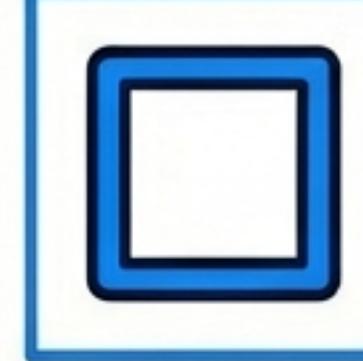
Spaces are often included in the token itself.

# The Final Translation: Text to Numbers

What the model actually sees.

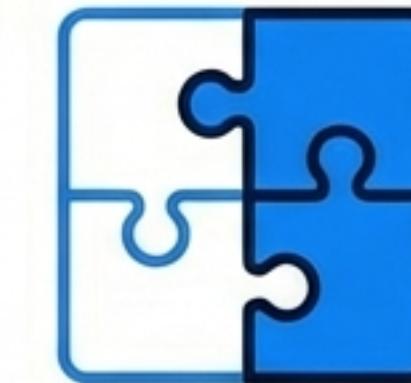


# Why Subword Tokenization Wins



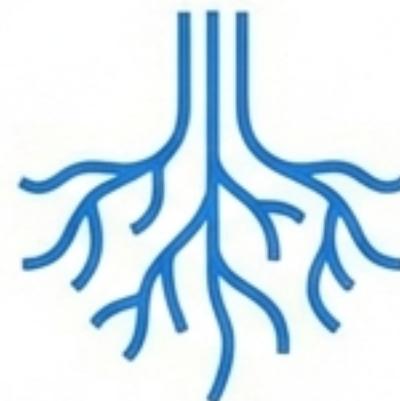
## Smaller Vocabulary

Efficient storage (30k-50k tokens vs millions).



## Handles Unknowns

Constructs new words from known parts.



## Preserves Structure

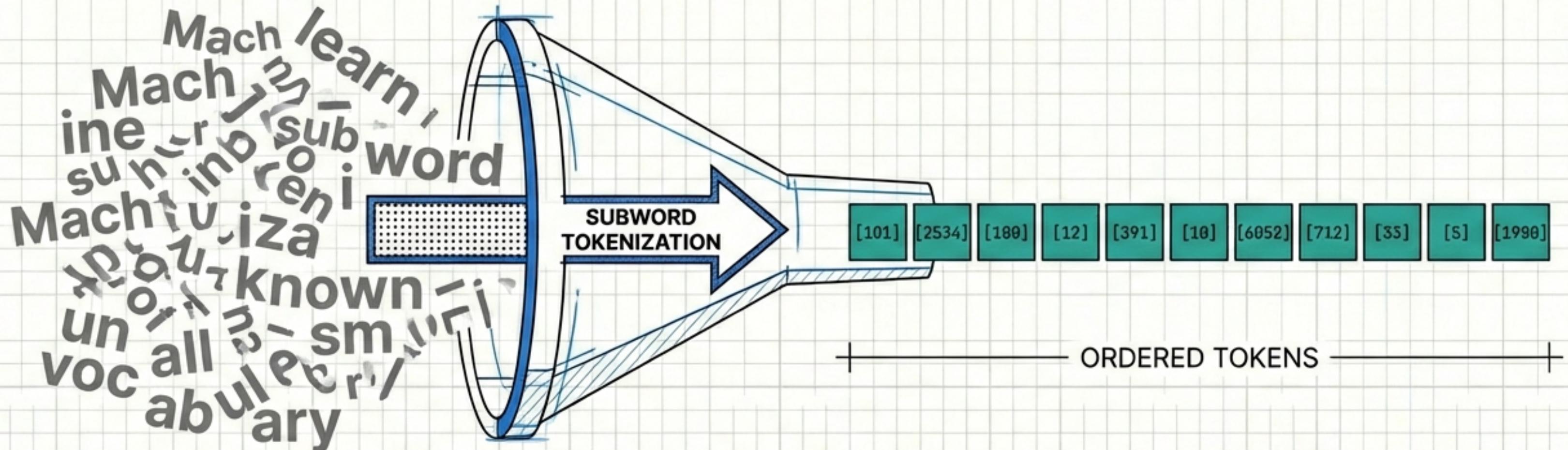
Keeps linguistic meaning (prefixes/suffixes).



## Multilingual

Universal standard for modern LLMs.

# The Takeaway



“ Subword tokenization balances efficiency with meaning. It splits words into frequent pieces so the model can understand the unknown while keeping the vocabulary small. ”