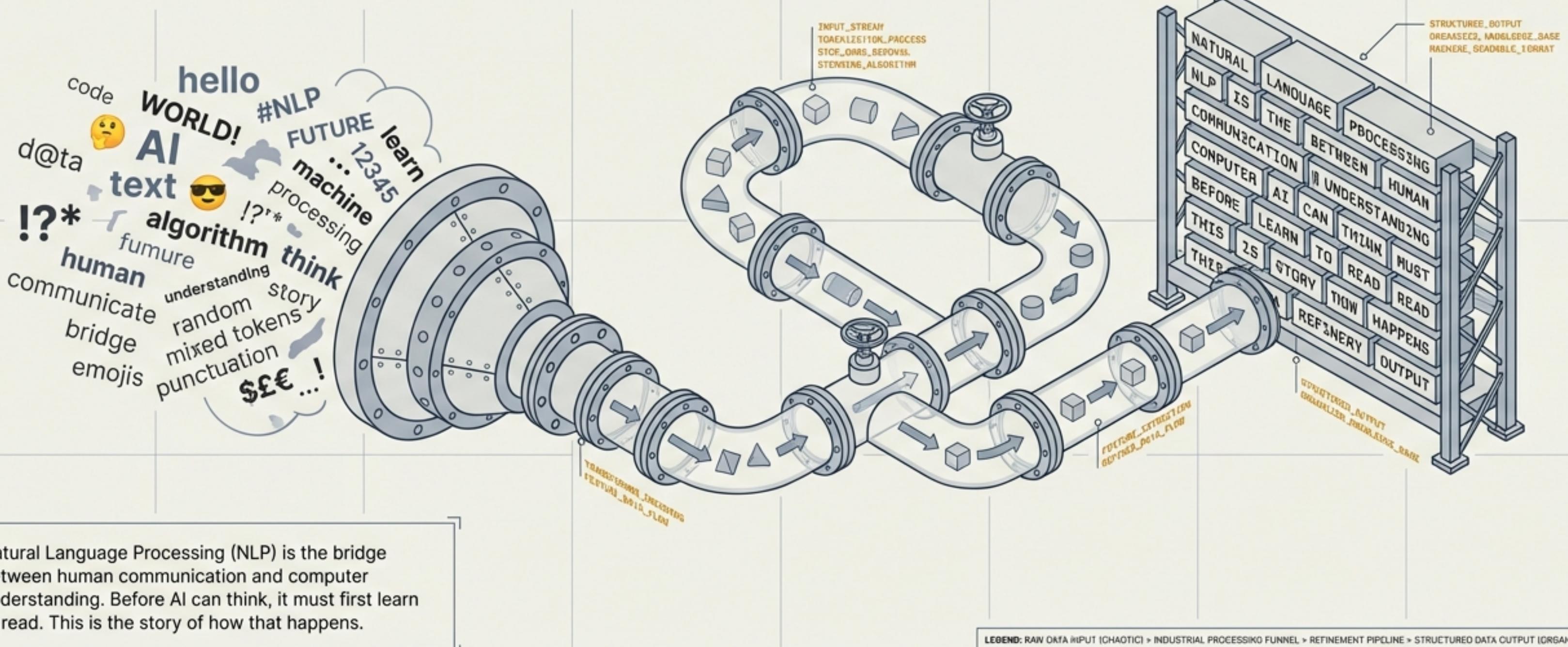
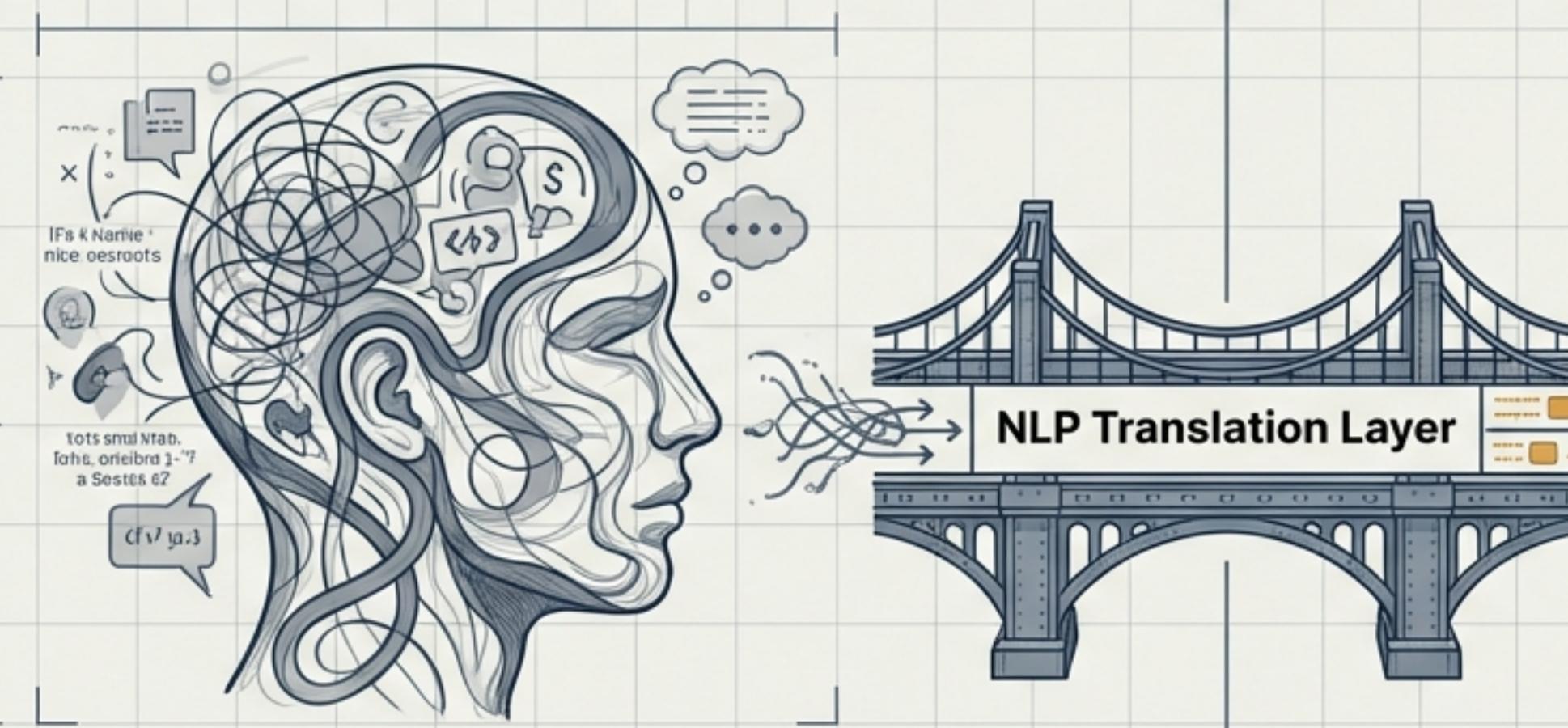


How Machines Read: The Data Refinery

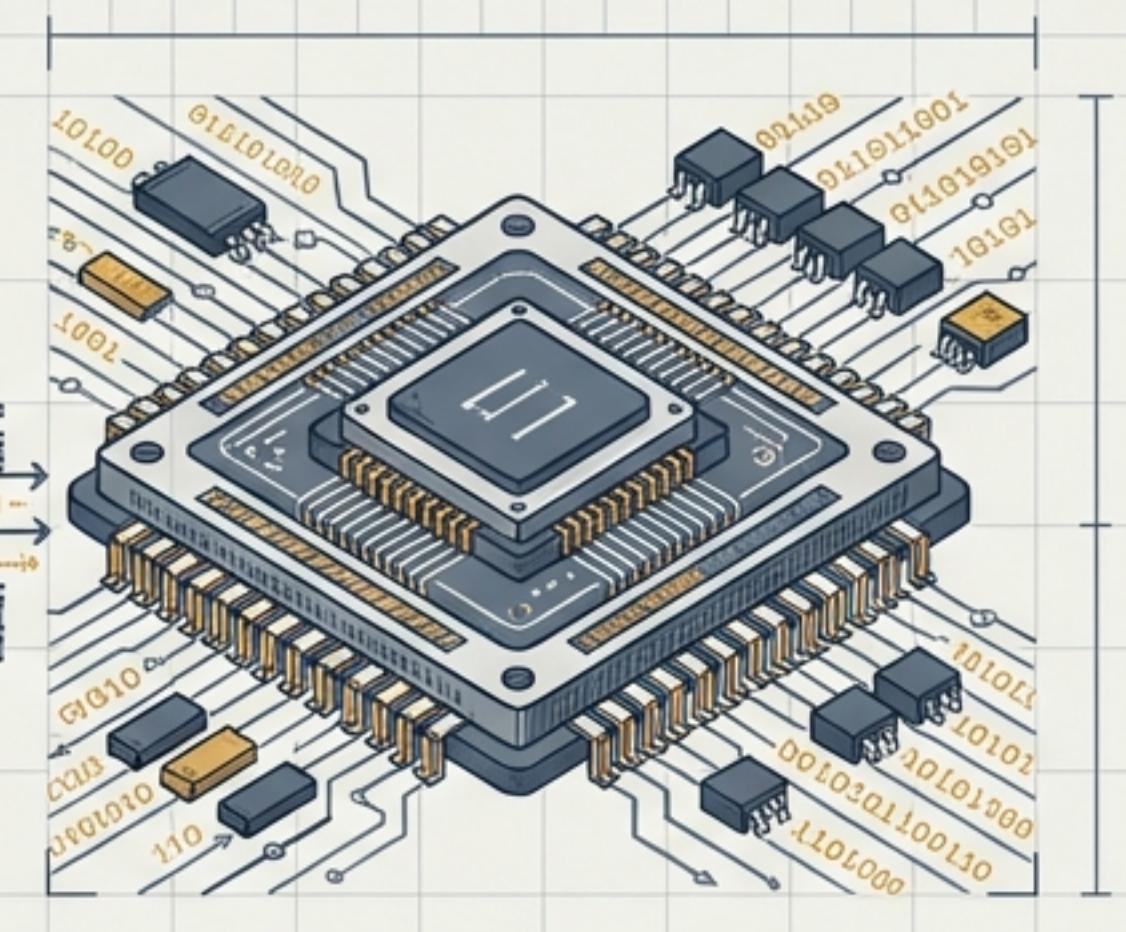
Turning human language into computer intelligence, one step at a time



The Language Barrier



Human: Ambiguity, Idioms, Context.



Machine: Logic, Math, Binary.

Natural Language Understanding (NLU)

Reading and interpreting.

E.g., Sentiment Analysis ("Is this email angry?").

Natural Language Generation (NLG)

Writing and responding.

E.g., ChatGPT (“Write a poem about cats”).

NLP enables machines to comprehend, generate, and manipulate human language.

You Have Already Used NLP Today



Spam Filters

"Is this junk?"
(Text Classification).

INBOX -> MODEL -> JUNK/NOT_JUNK



Voice Assistants

"Siri, set an alarm."
(Speech-to-Text).

AUDIO -> DECODE -> 'SET ALARM'



Translation

"Coffee in Paris?"
(Machine Translation).

EN -> MODEL -> FR -> 'CAFÉ À PARIS?'



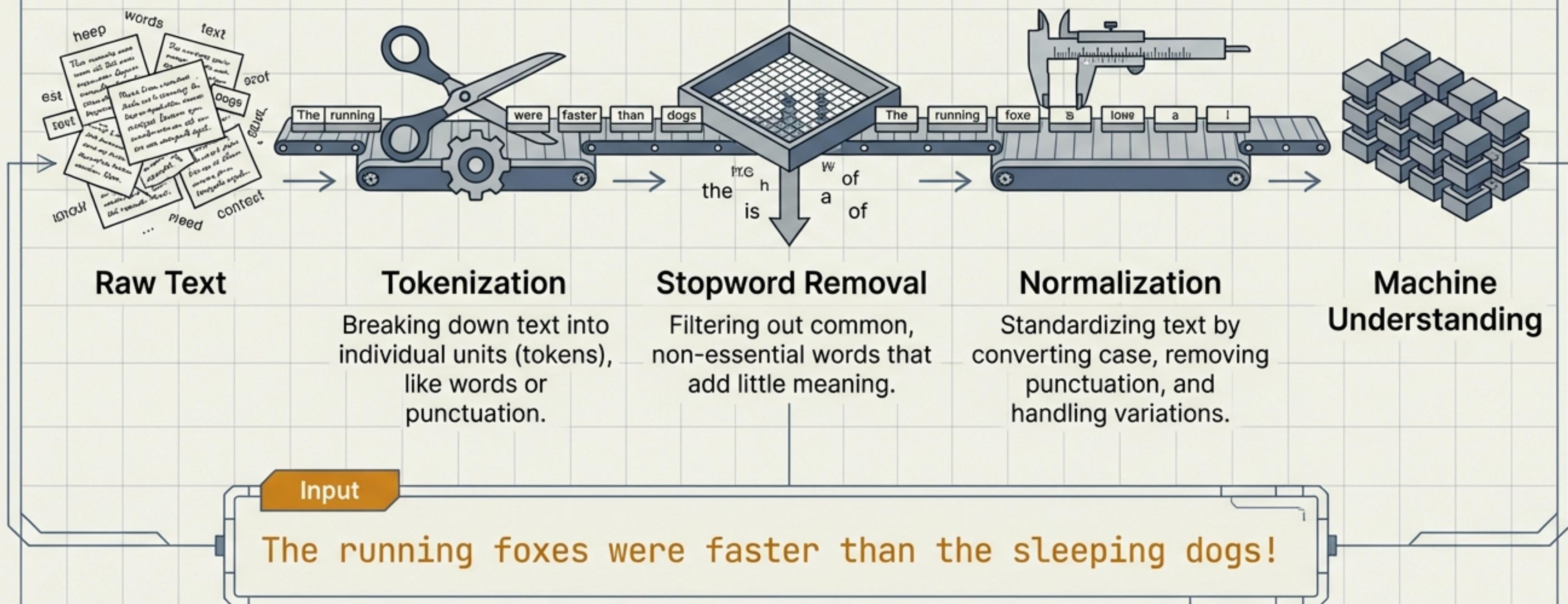
Autocorrect

'Teh' -> 'The'
(Spell Check).

'Teh' -> SPELL_CHECK -> 'The'

The Insight: NLP turns unstructured text (which computers hate) into structured insights (which computers need).

The Processing Pipeline



Step 1: Tokenization

Breaking the stream into manageable chunks.

Before

The running foxes were faster than the sleeping dogs!

After



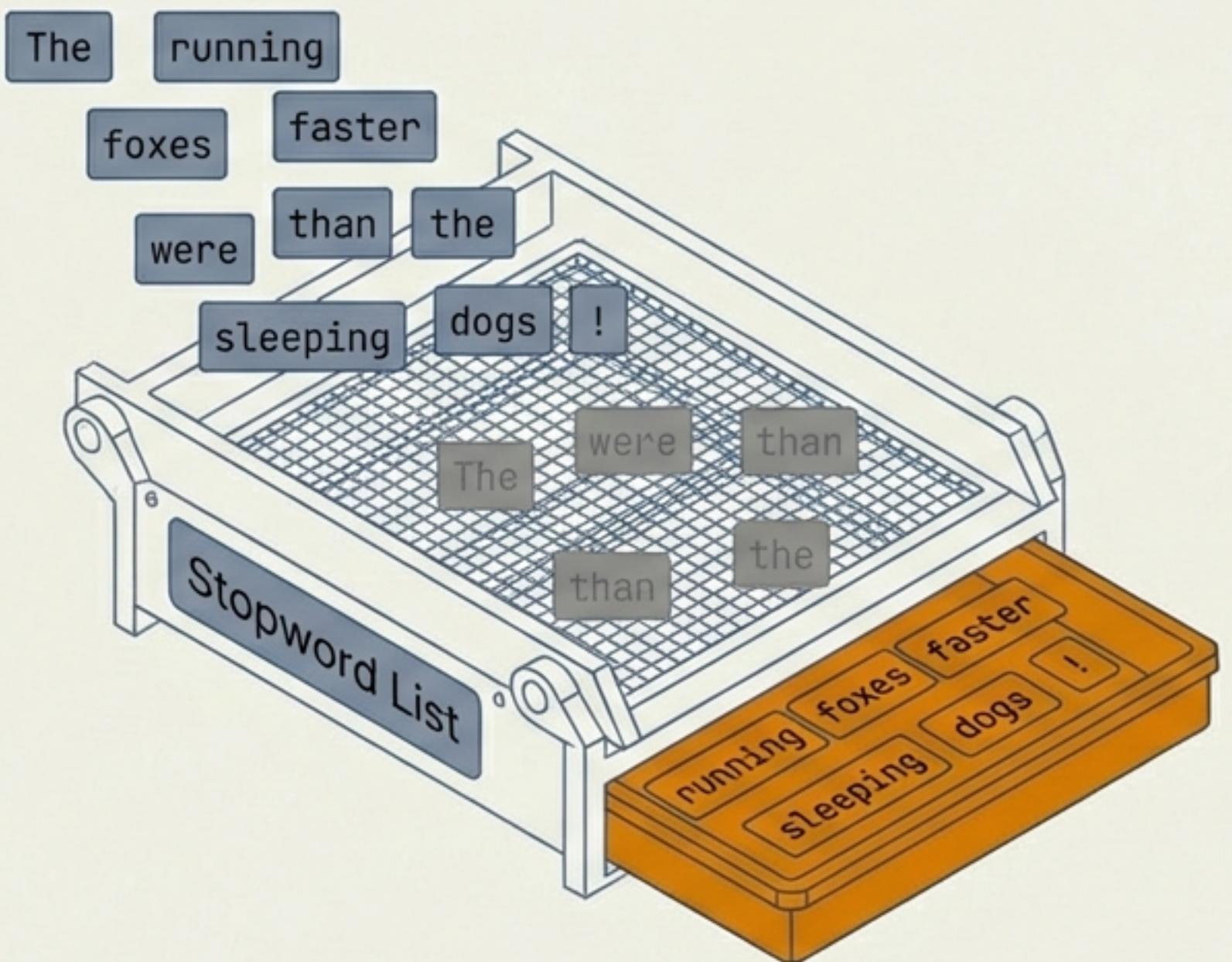
Tokenization is the foundational step of breaking a stream of text into smaller units called tokens.

It's Not Just About Whole Words

Word Tokenization	[Unhappiness]	Standard approach. Large vocabulary size.
Sentence Tokenization	A 3D orange parallelogram labeled "PARAGRAPH" points to two adjacent orange boxes. The top box contains the text "The first sentence runs here." and the bottom box contains the text "The second sentence follows closely."	Splits paragraphs. Good for summarization.
Character Tokenization	[U] [n] [h] [a] [p] [p] [i] [n] [e] [s] [s]	Granular. Handles typos, but computationally expensive.
Subword Tokenization (BPE)	[Un] [happi] [ness]	Used by GPT-4. Breaks unknown words into meaningful roots.

Step 2: Stopword Removal

Filtering the noise to focus on signal.



Stopwords are high-frequency words (is, the, at) that add structure but little unique meaning. Removing them reduces data size and highlights keywords.

The 'Stopword' Trap

⚠ CAUTION

To be or not to be



If we remove stopwords, this famous phrase disappears entirely.

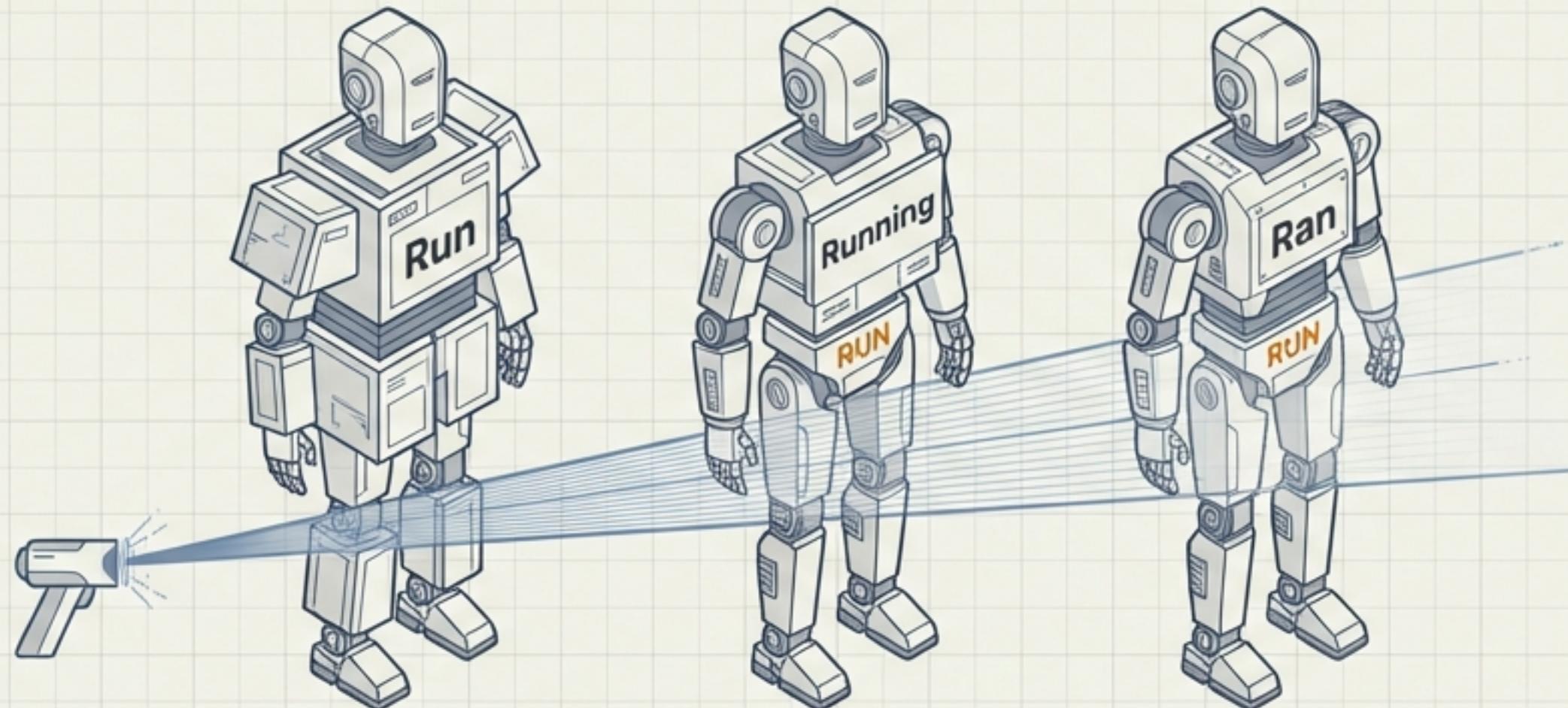


In advanced AI, stopwords provide critical context. 'From' and 'To' define the destination.

Modern LLMs (like GPT) often keep stopwords to preserve full semantic meaning.

The Problem of Disguises

Normalization: Teaching the machine that “Run” and “Running” are the same.

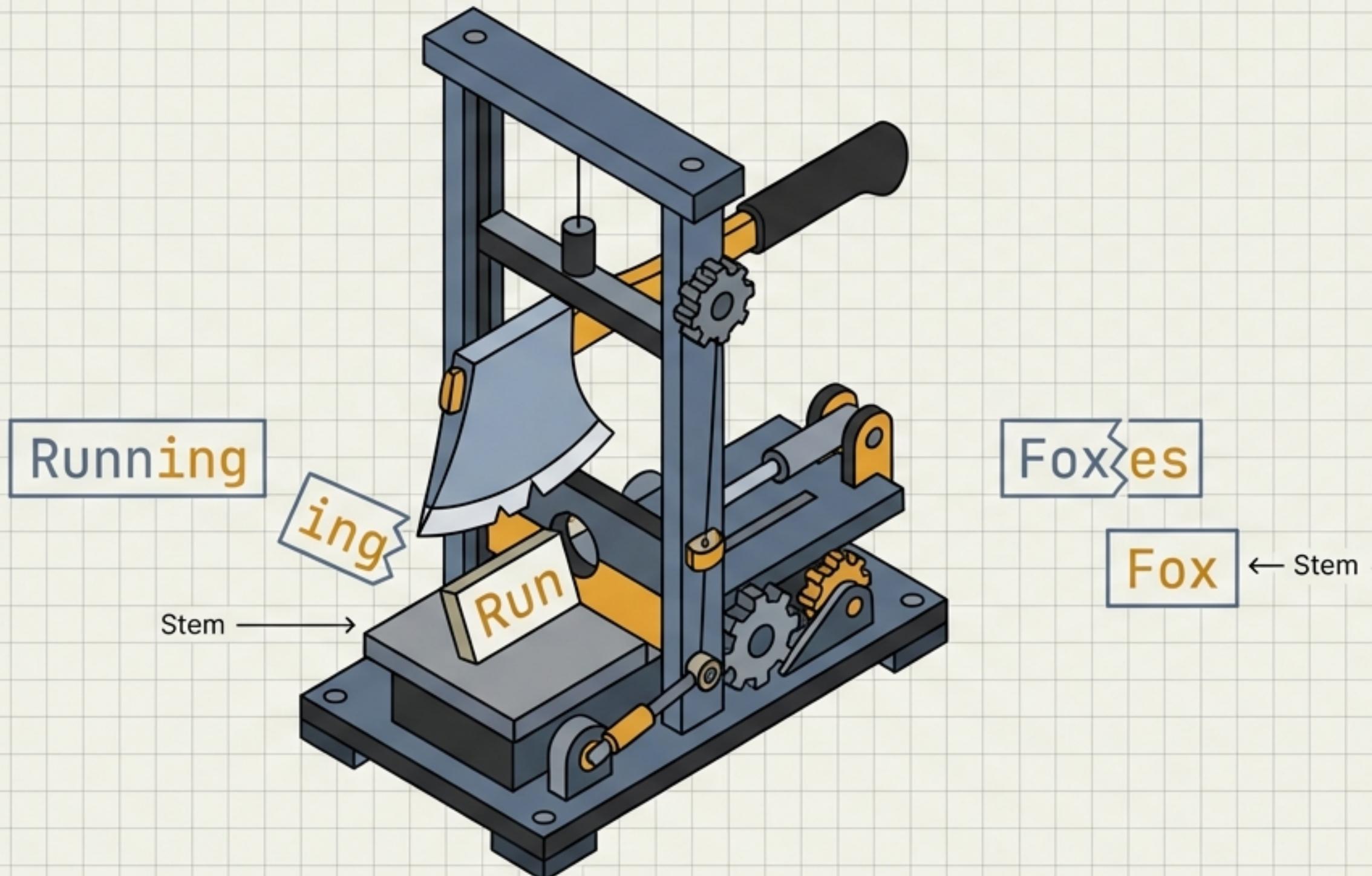


Computers see these as three mathematically different words.

Normalization strips the disguises.



Method A: Stemming (The Hatchet)



Stemming is a rule-based approach that aggressively chops off suffixes.

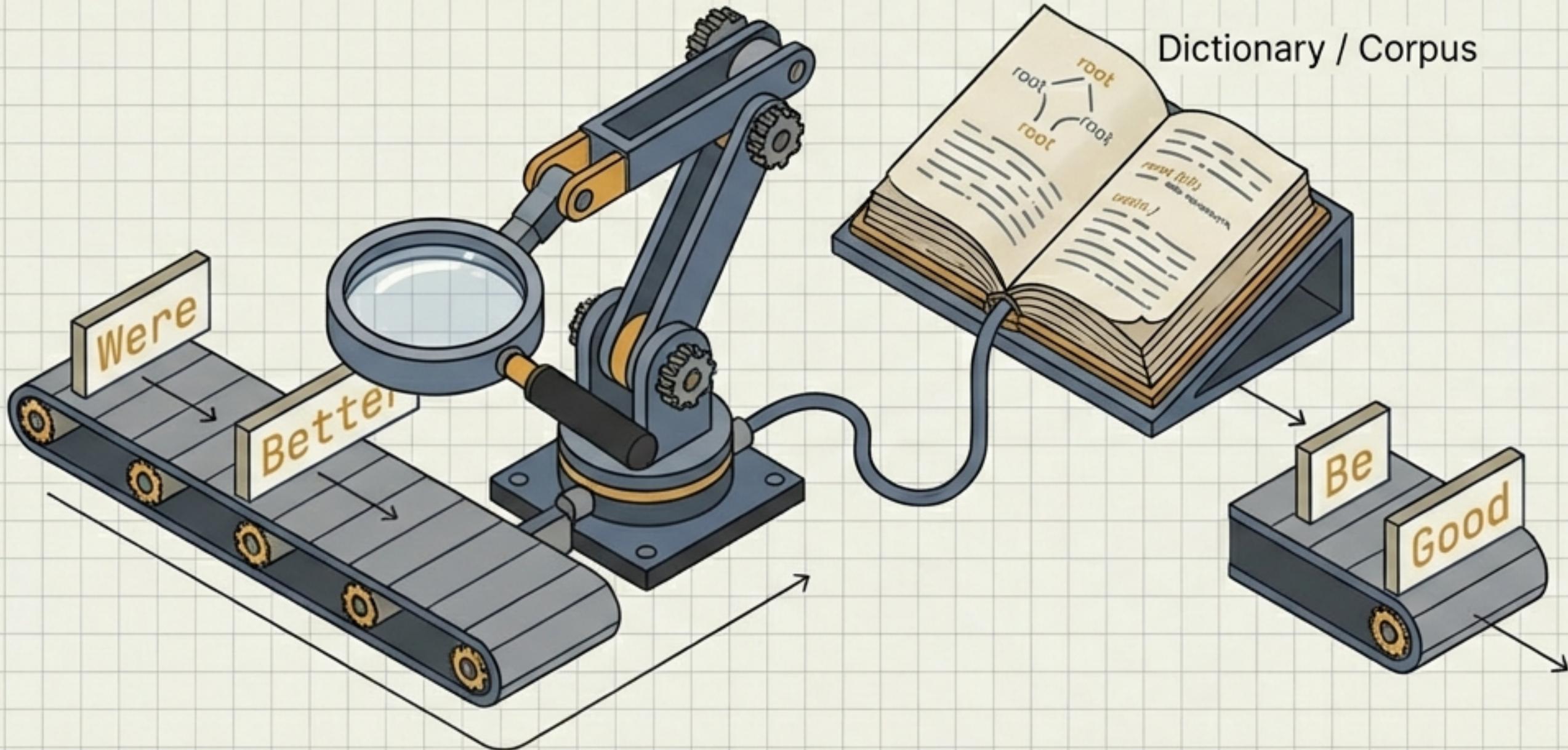
The Flaw: Over-stemming.

- Example: "Ponies" -> chopped to -> "Poni" (Not a real word).
- Example: "Universal" -> chopped to -> "Univers".

Pros

- ✓ Extremely Fast.

Method B: Lemmatization (The Librarian)



Lemmatization looks up the word in a dictionary to find its true linguistic root (**Lemma**). It understands that 'Better' is the comparative form of 'Good'.

✓ Pros:

Highly Accurate.
Always produces
real words.

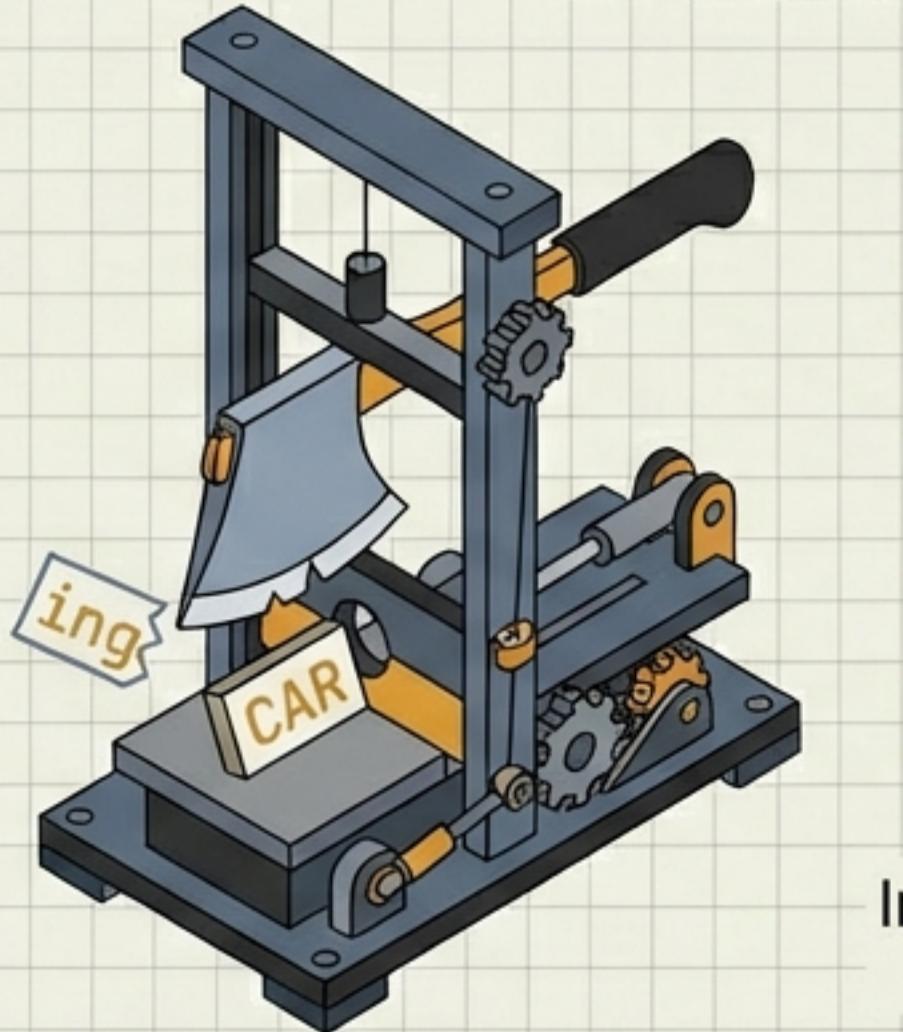
✗ Cons:

Computationally
Slower.

Showdown: Stemming vs. Lemmatization

Stemming

The Hatchet



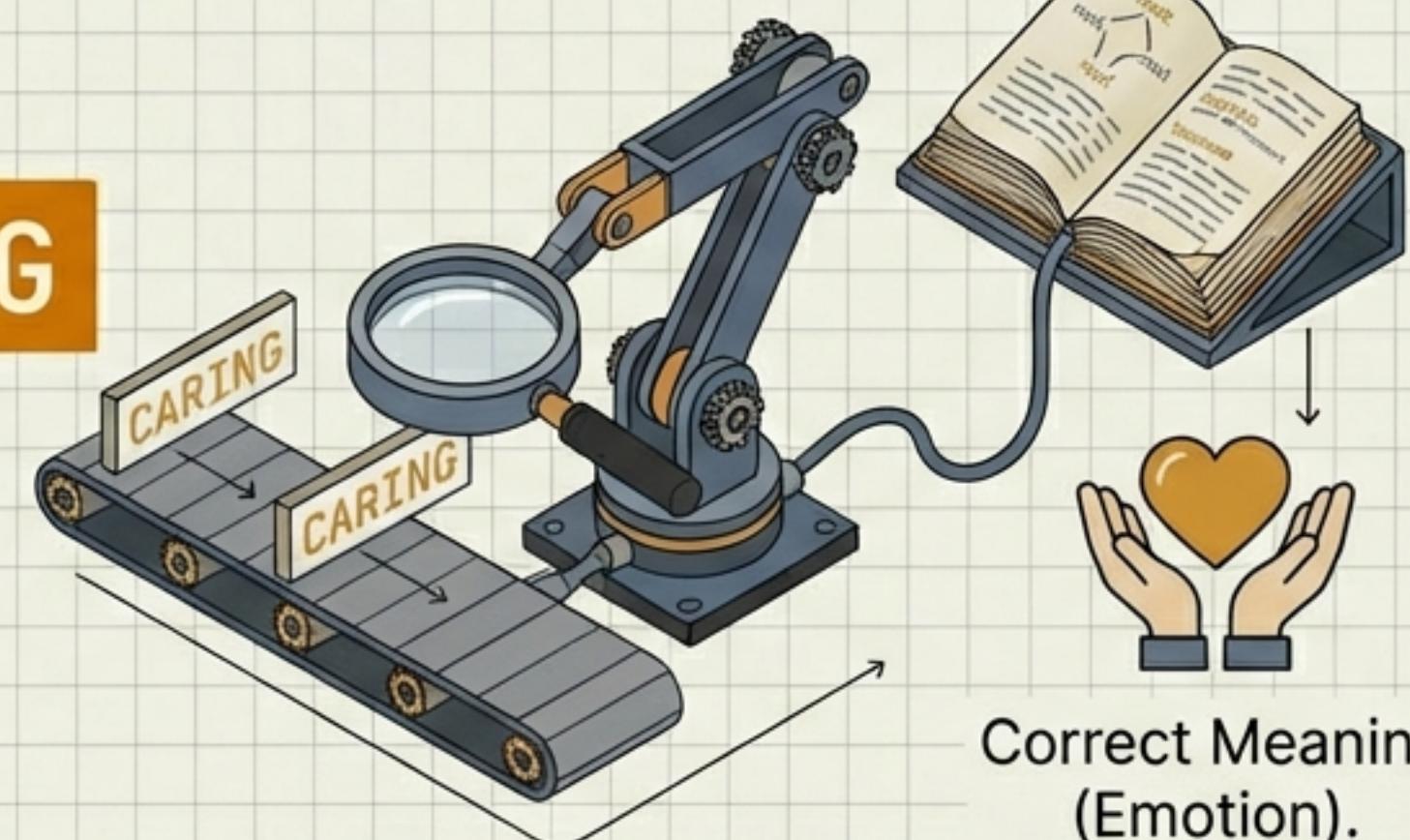
Incorrect Meaning
(Vehicle).

CARING



Lemmatization

The Librarian



Correct Meaning
(Emotion).

Use **Stemming** for speed (Search Engines). Use **Lemmatization** for accuracy (Chatbots/Translation).

The Full Transformation

Step 1. Original

The running foxes were faster than the sleeping dogs!

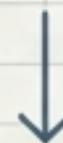
Step 2. Tokenized

The running foxes were faster than the sleeping dogs !!



Step 3. Stopwords Removed

The running foxes faster than sleeping dogs !



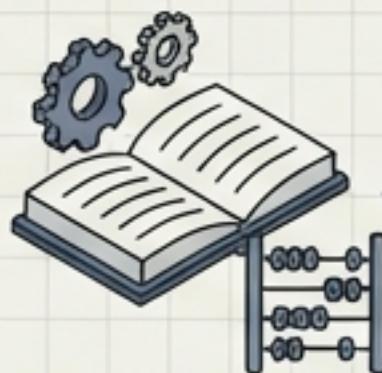
Step 4. Lemmatized (Final Output)

run fox fast sleep dog !



The machine now receives a streamlined list of core concepts, ready for analysis.

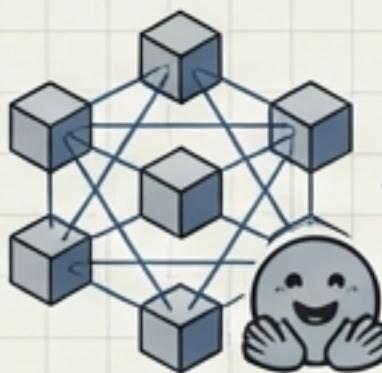
The Toolkit



NLTK: Academic Standard.
Great for learning.



spaCy: Industrial Strength.
Fast and production-ready.



Hugging Face: Transformer Era.
The modern standard for LLMs.

A screenshot of a Python code editor window titled "Python - lemmatization.py". The code is as follows:

```
1 import spacy
2 nlp = spacy.load('en_core_web_sm')
3 doc = nlp('The running foxes...')
4 print([token.lemma_ for token in doc])
```

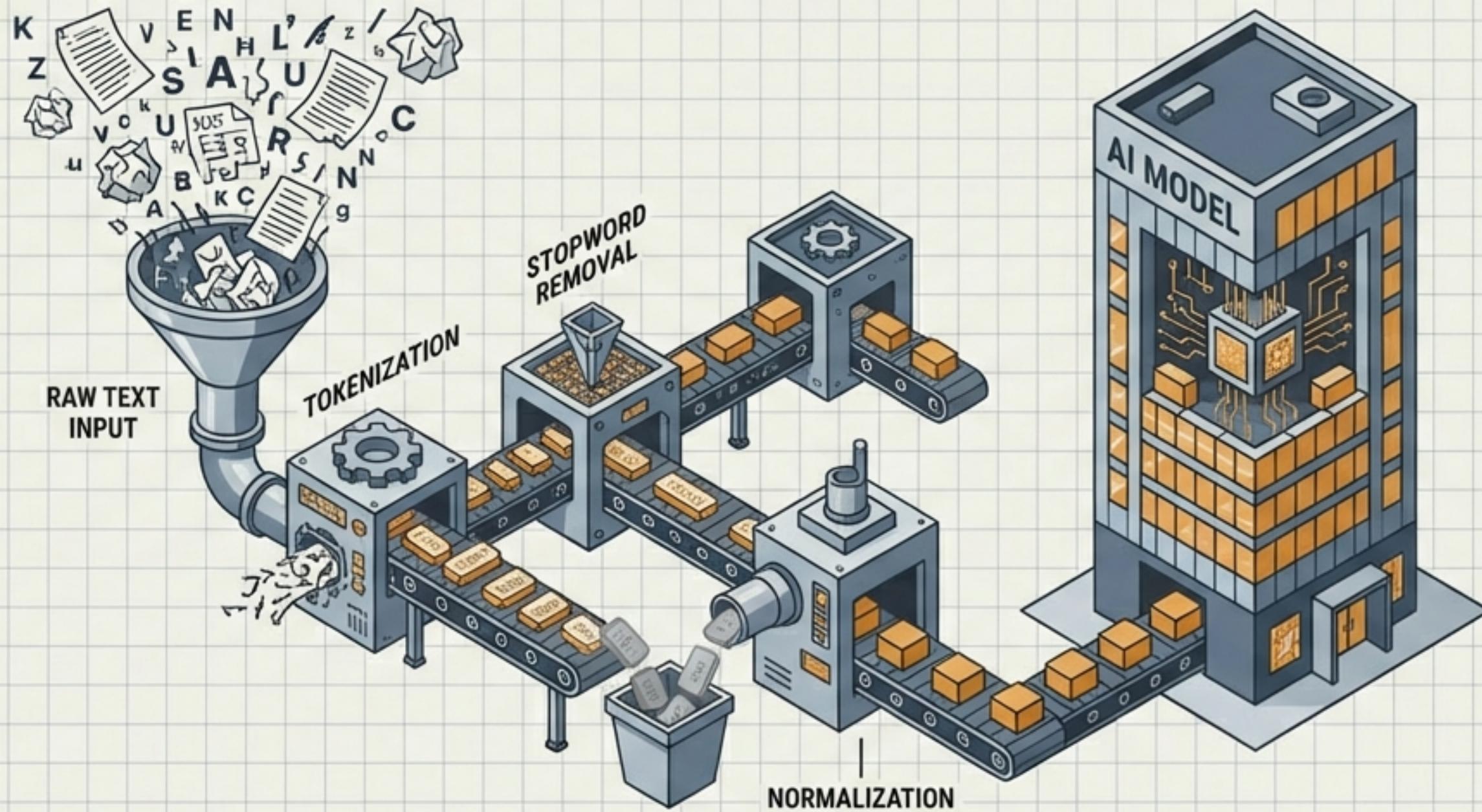
A screenshot of a terminal window titled "Output" showing the result of the executed code:

```
['the', 'run', 'fox', '...']
```

From Chaos to Order

Summary

1. **Tokenization** breaks the bricks.
2. **Stopwords** clear the debris.
3. **Normalization** shapes the bricks for building.



Language is messy. NLP Preprocessing is the cleanup crew that organizes this mess so AI can begin to understand us.