Before Transformers, models like RNNs (Recurrent Neural Networks) and LSTMs (Long Short-Term Memory) were used to process sequences (like sentences) 1 step at a time — sequentially. Word 1  $\rightarrow$  Word 2  $\rightarrow$  Word 3  $\rightarrow$  ...  $\rightarrow$  Word N Each word depends on the one before it. That meant: Transformers introduced by Google This is slow and hard to parallelize, especially for long texts ("Attention is All You Need") Each word in a sentence looks - the atomic bomb of GenAI. at every other word all at once. Instead of recurrence, they use self-attention: The model figures out which words to pay attention to, no matter their position The mother hit the daughter because she was drunk Bidirectional understanding of language. BERT (2018, Google): Powers Google Search and beyond. Shockingly good at generating coherent text. GPT-2 (2019, OpenAI): OpenAl withholds full release initially due to misuse concerns. Language Models Unified NLP framework: every NLP task Take Off becomes a "text-in, text-out" problem. Transformers introduced by Google T5 (Text-to-Text Transfer Transformer, Google, 2020) ("Attention is All You Need) - the atomic bomb of GenAI. 175 billion parameters, trained on a massive corpus. 2020: OpenAl drops GPT-3 – shock and awe. Text-to-image synthesis (DALL·E), Raw audio generation (Jukebox), Code generation (Codex) Scale becomes the new weapon - billions of parameters, massive compute.

world war 2 (2017-2022) The transformer revolution