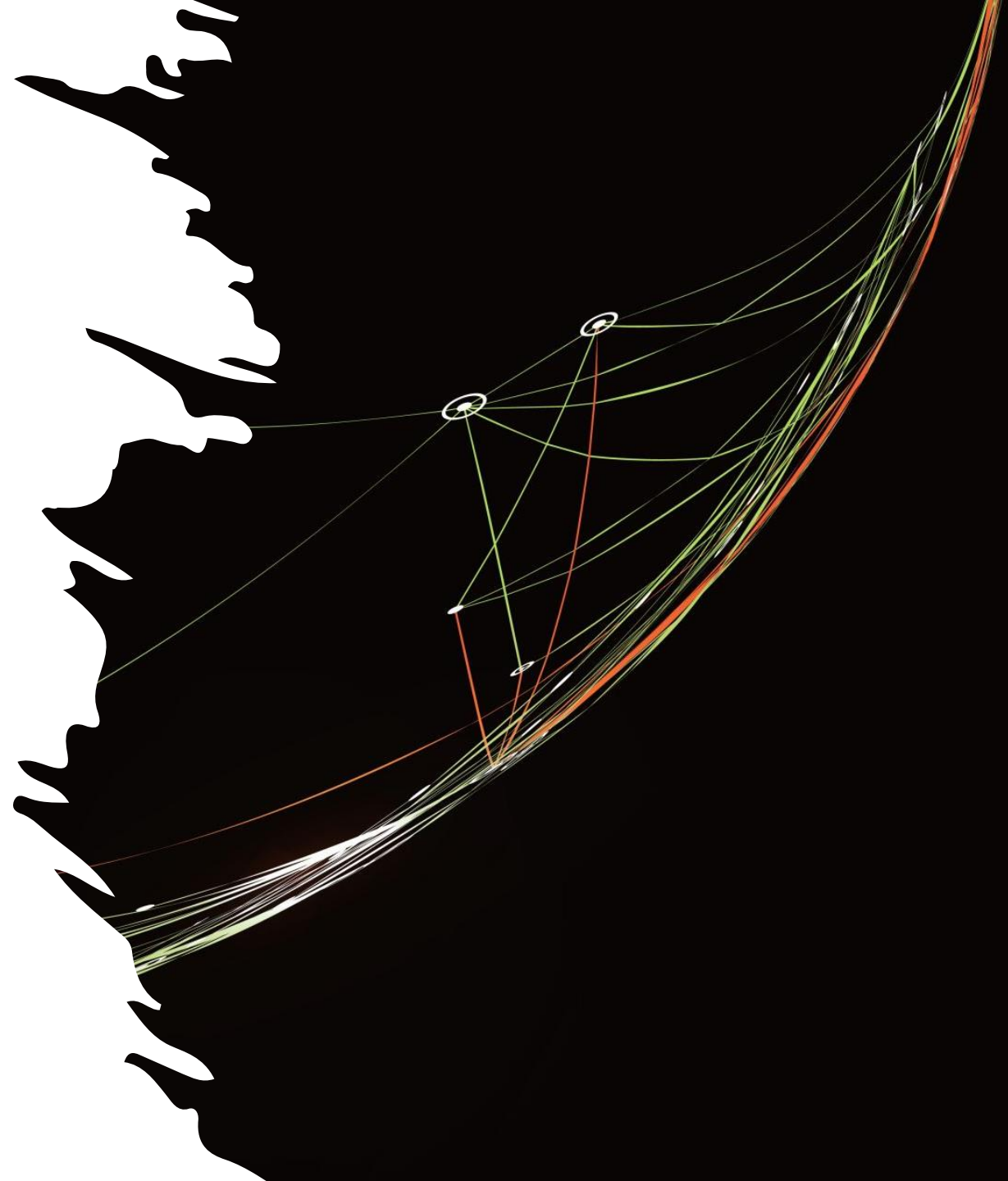
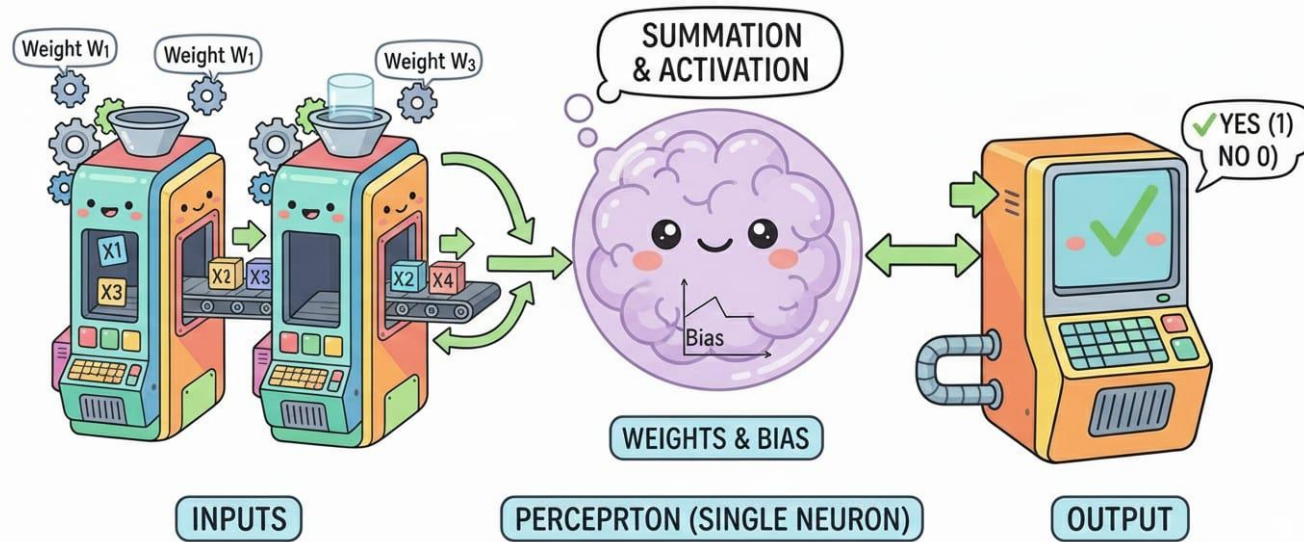


Un breve camino (lúdico) por las redes neuronales

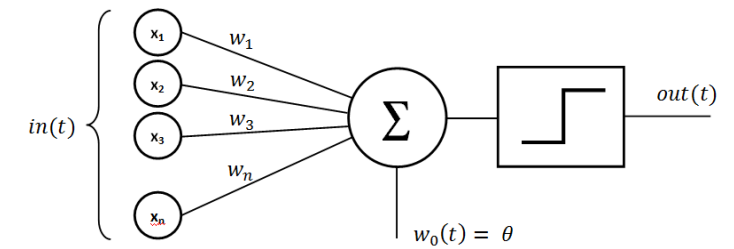
Erick Merino

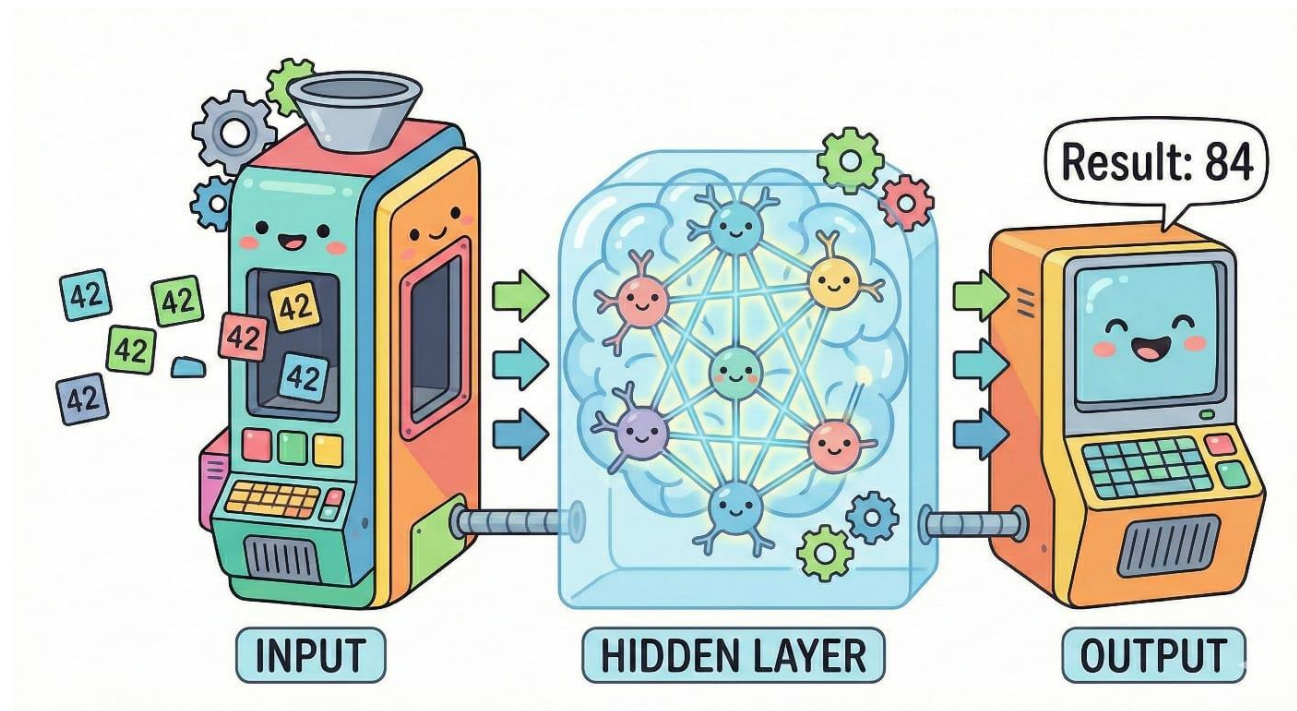


Perceptron

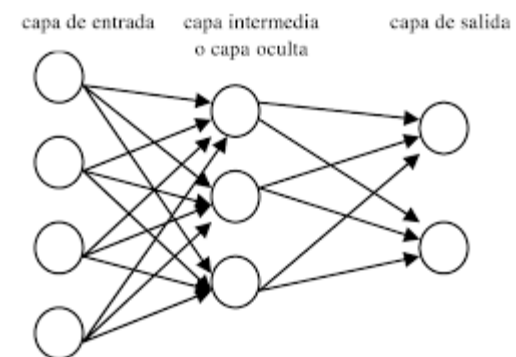


Todo empezó con una neurona

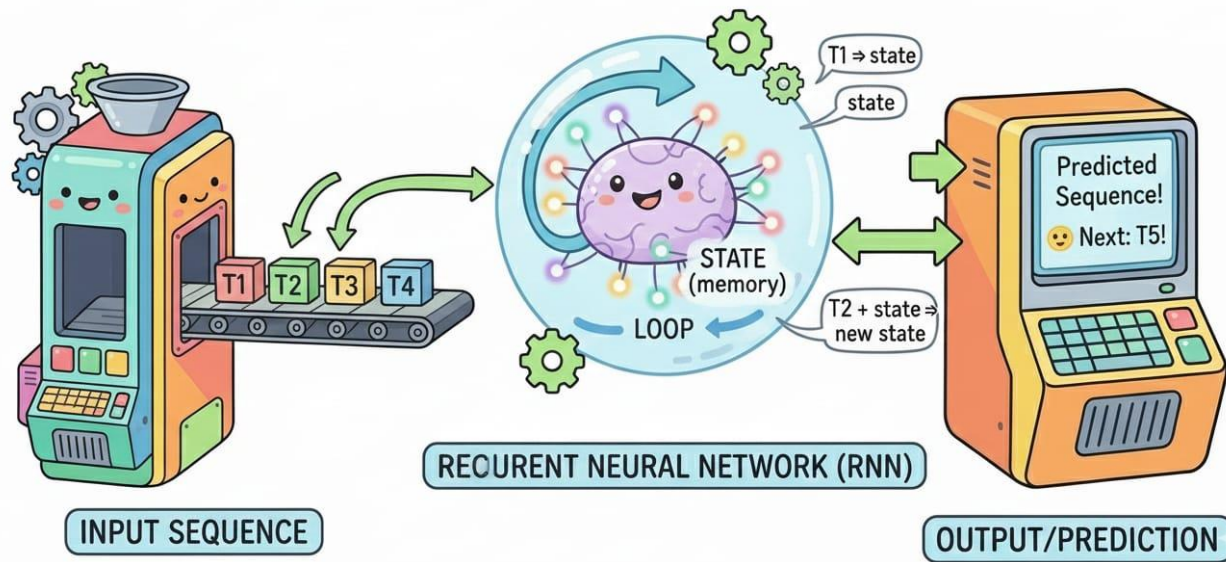




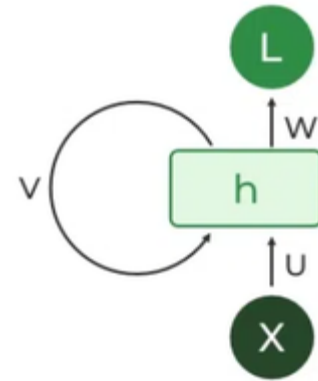
A la cual se le
sumaron más



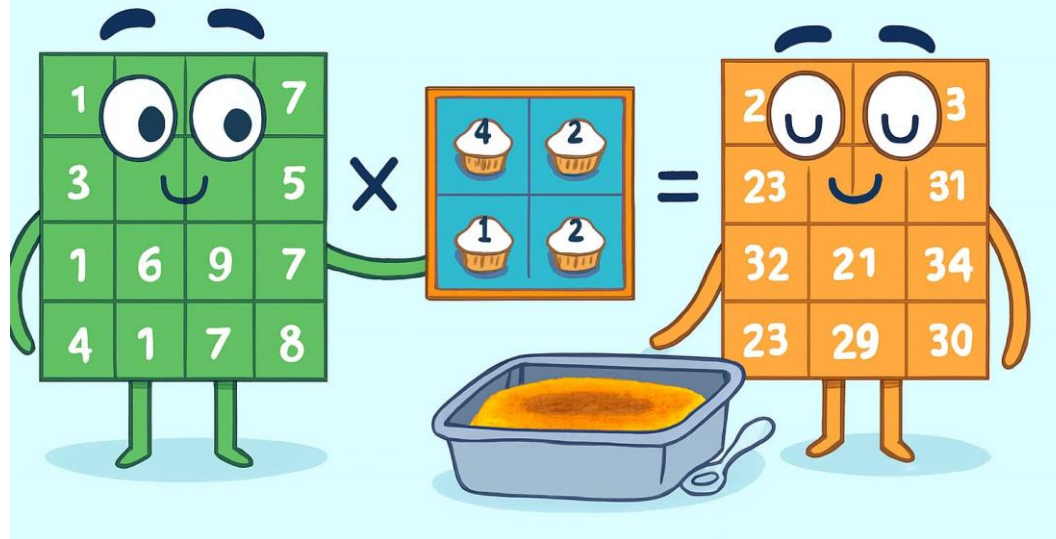
RNN: The Memory Machine!



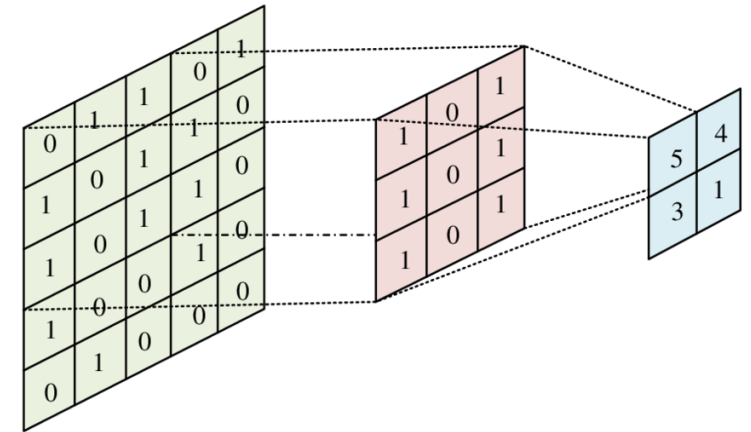
Se les dio memoria



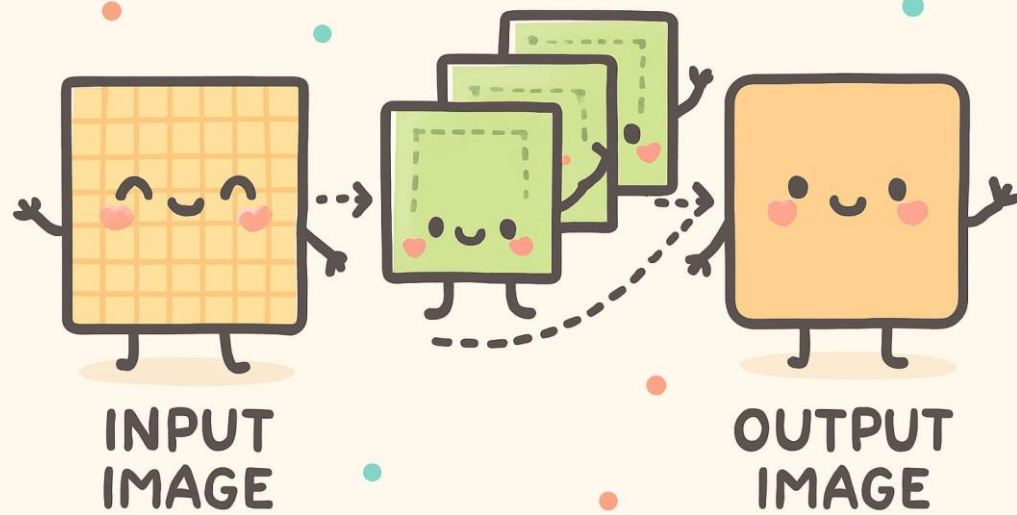
Convolution



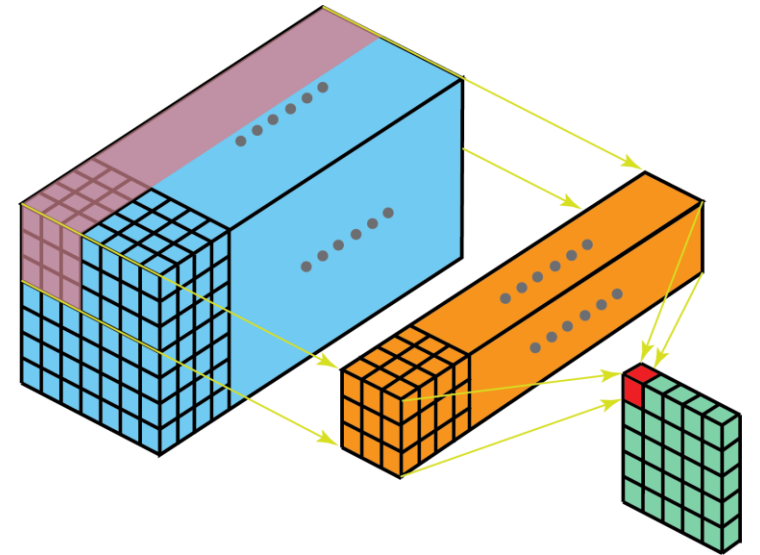
Se las ordenó y
multiplicó



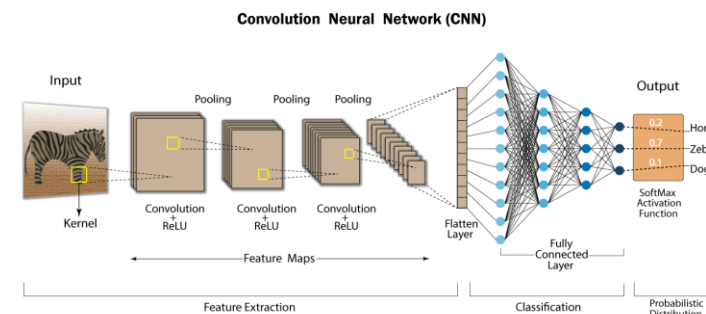
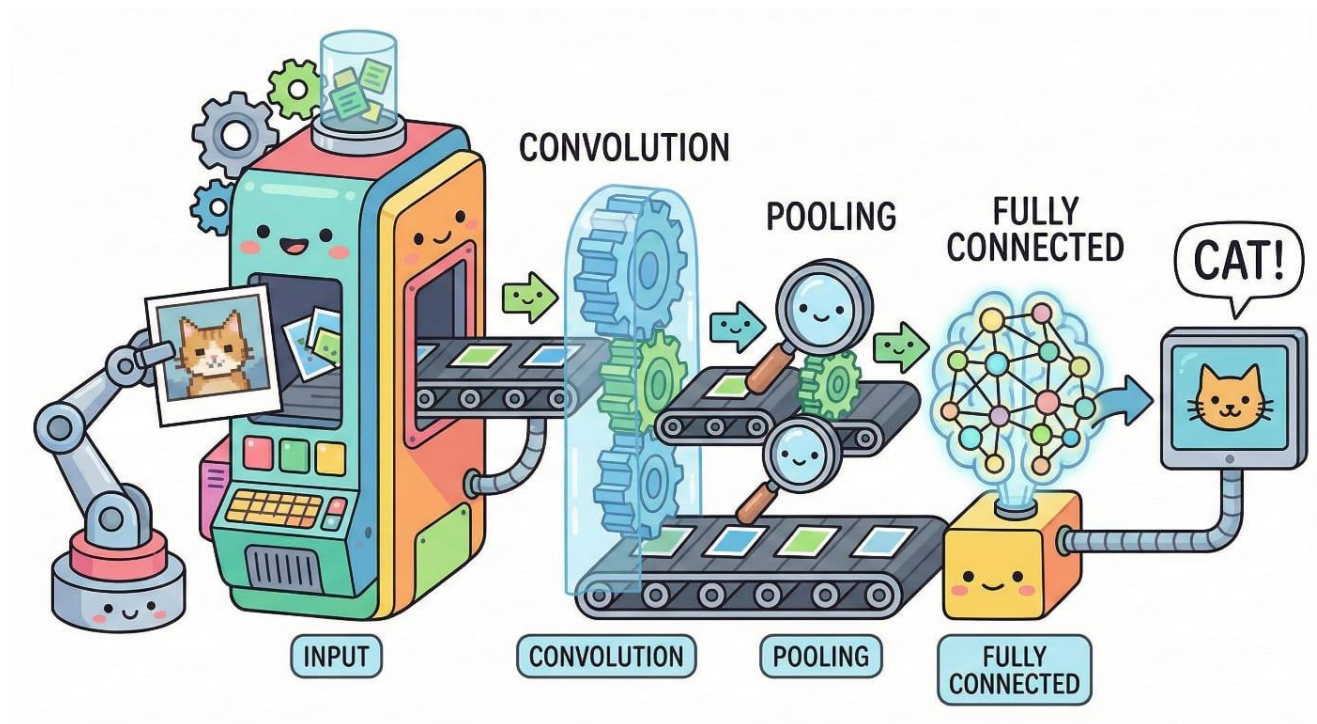
CONVOLUTION



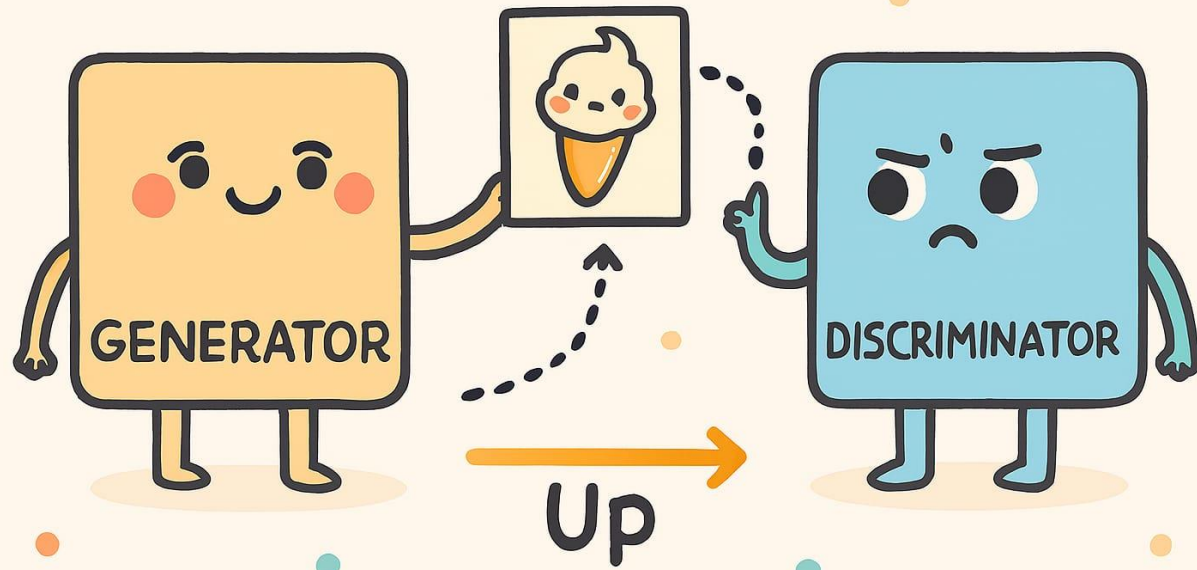
Múltiples veces



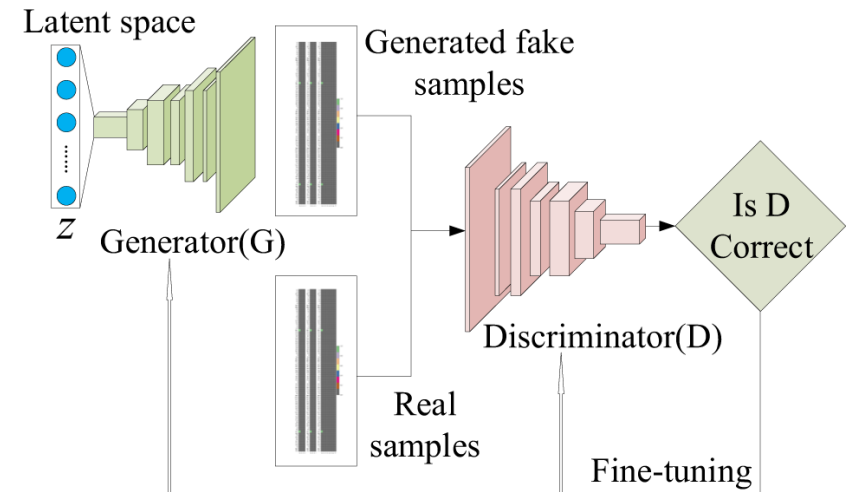
Se crearon complejas
máquinas de
neuronas



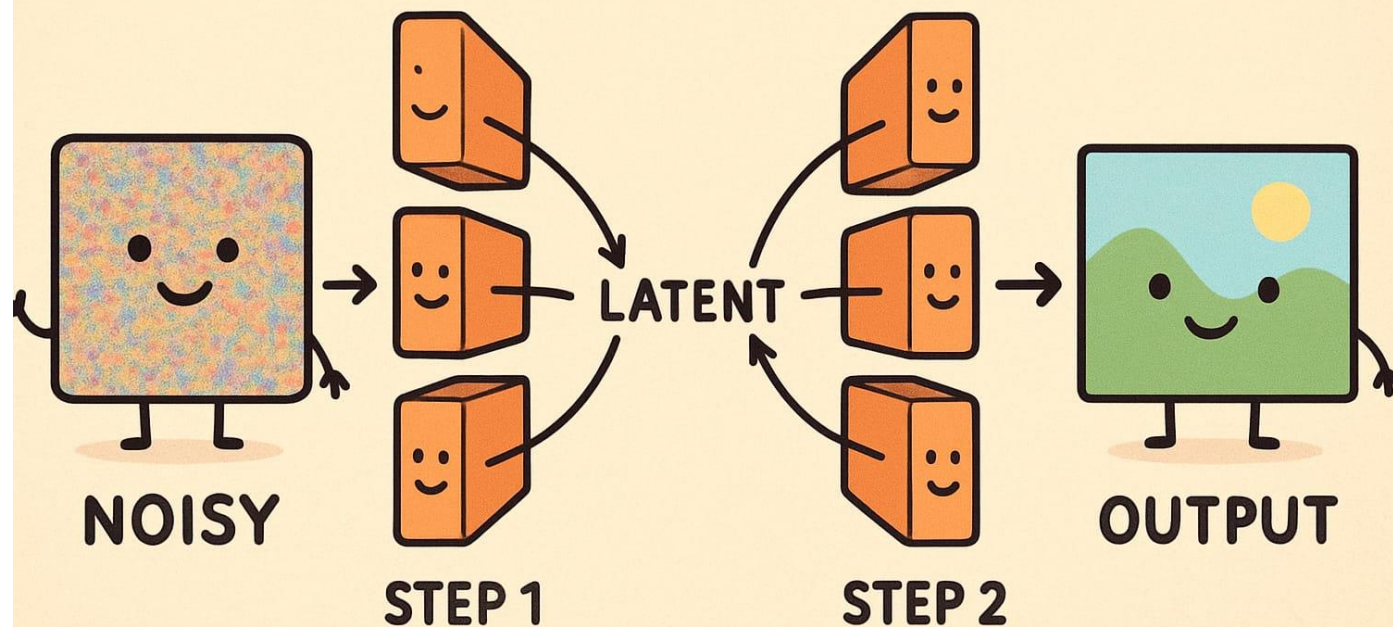
GAN



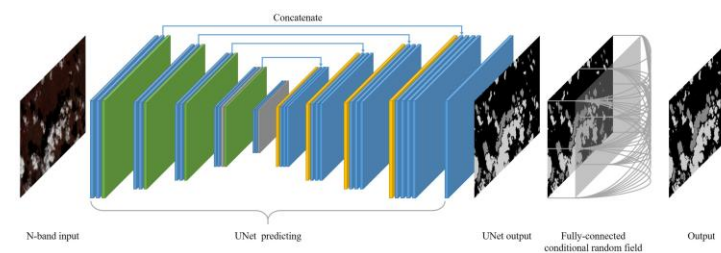
Se las hizo competir

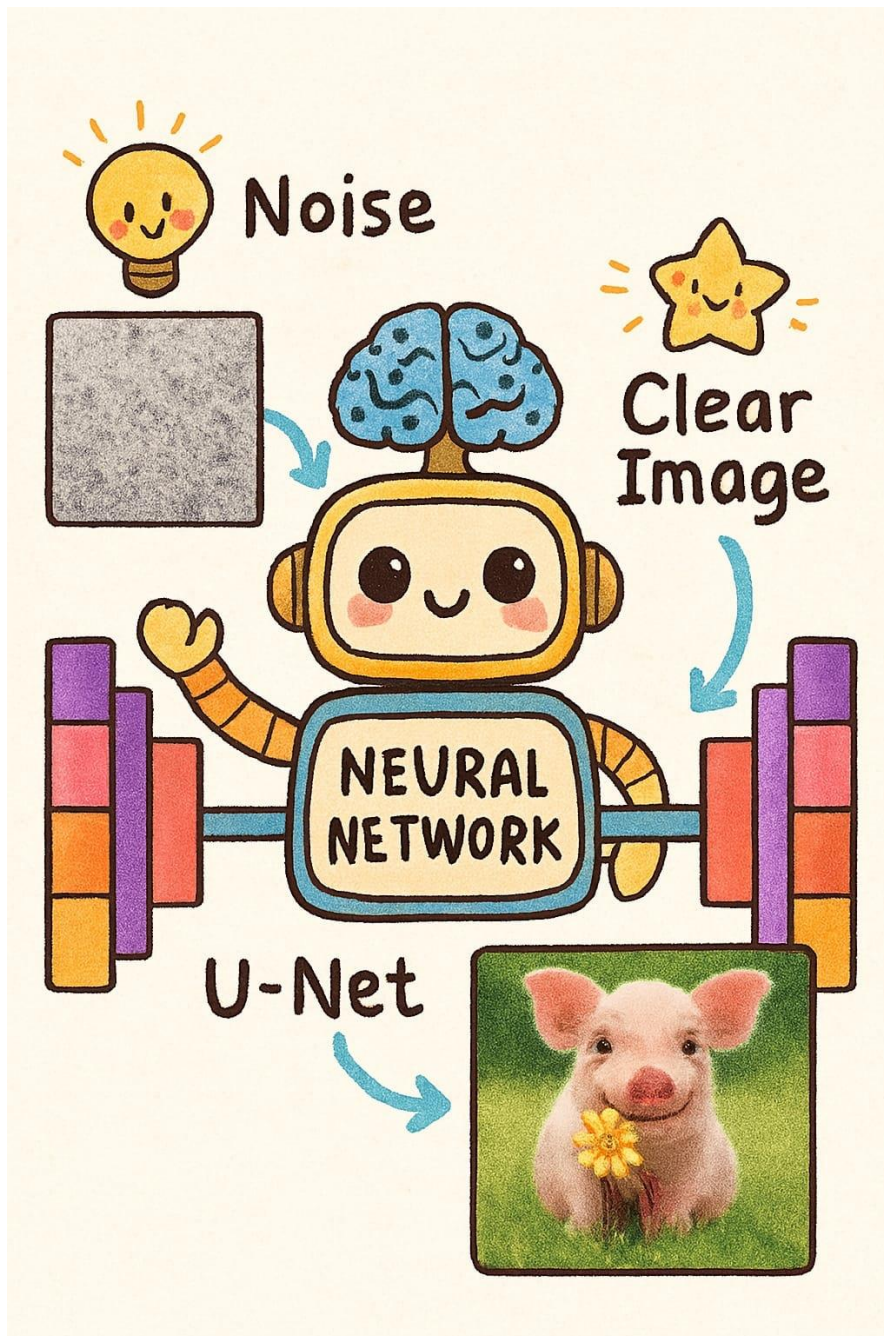


U-NET

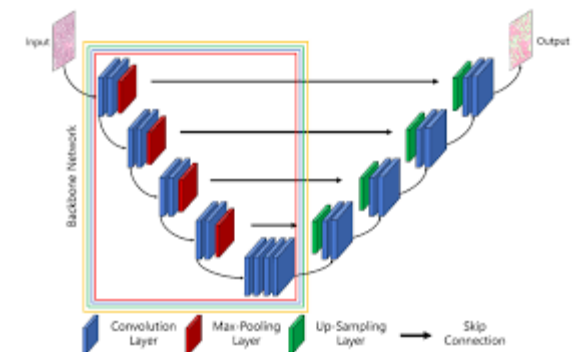


Se las
empequeñeció y
agrandó

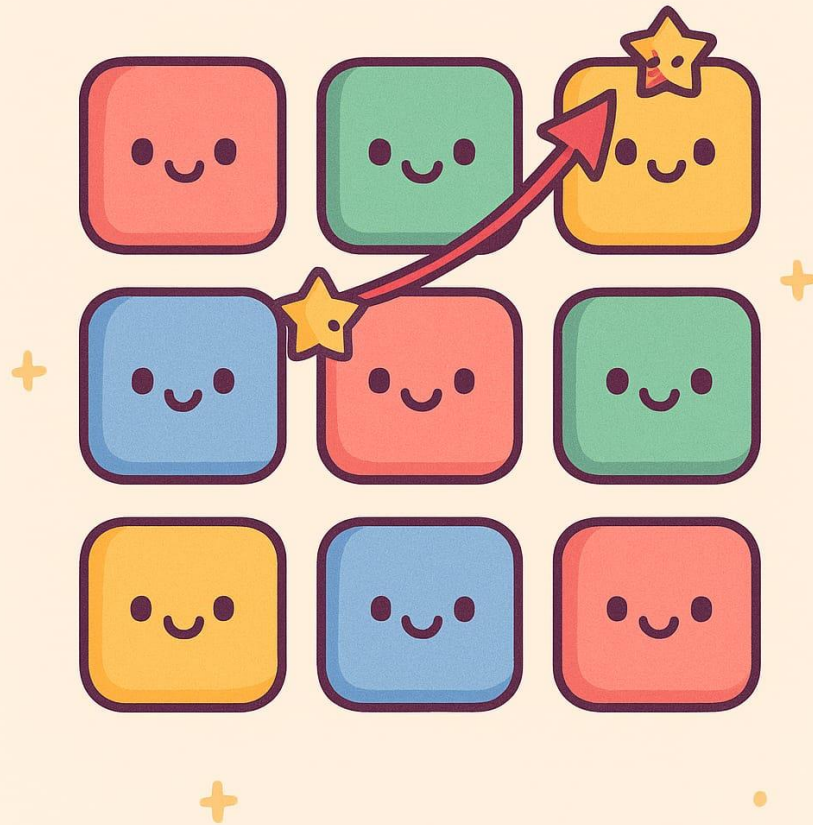




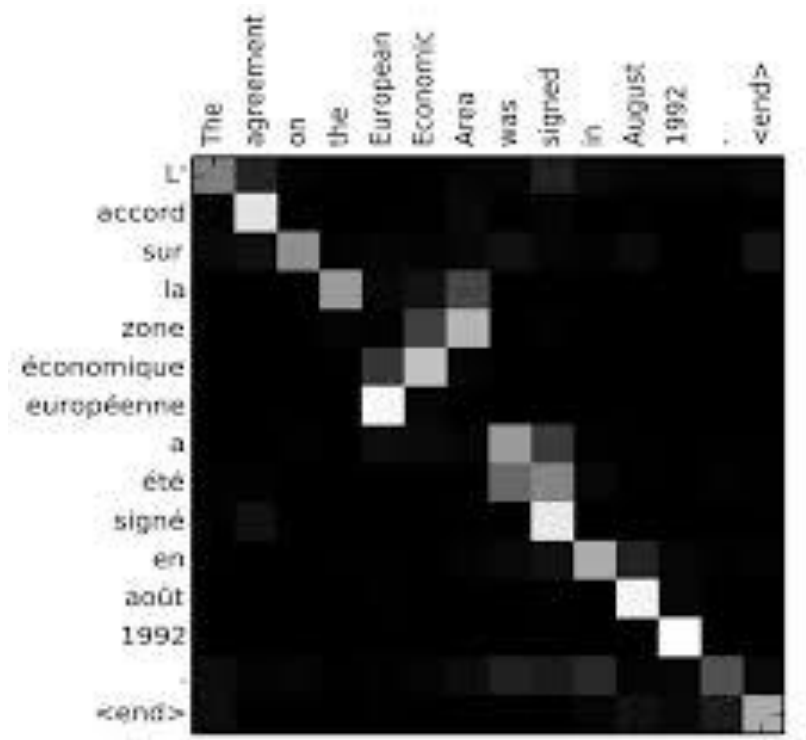
Fueron capaces de generar imágenes



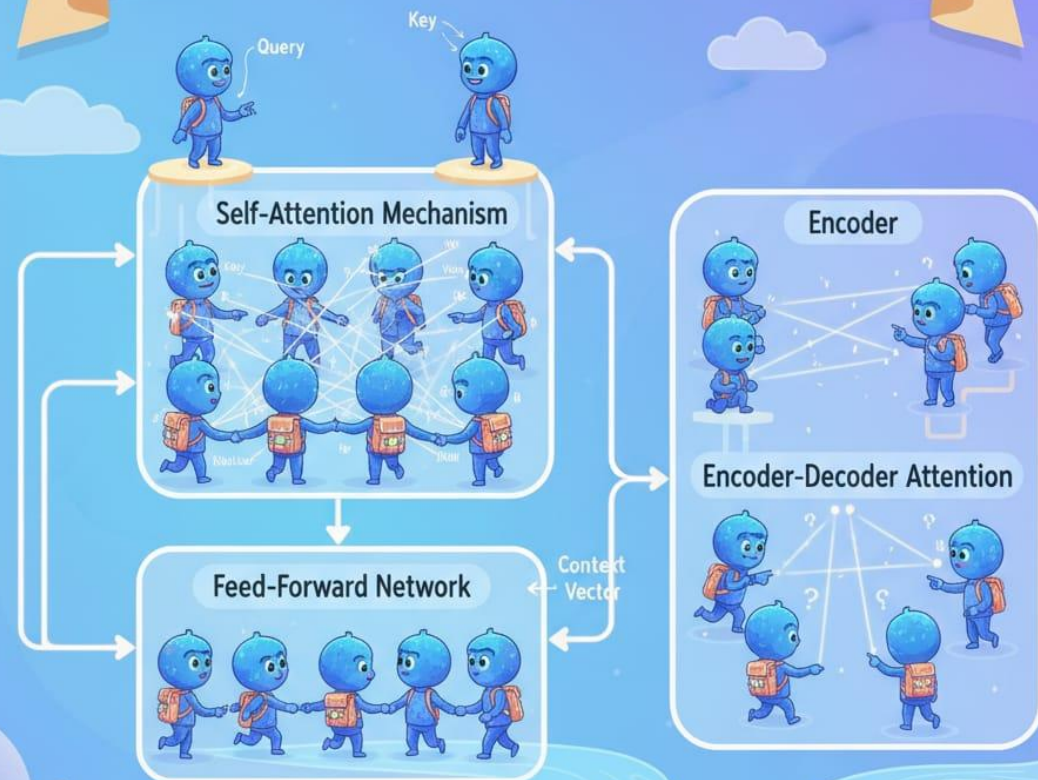
ATTENTION MATRIX



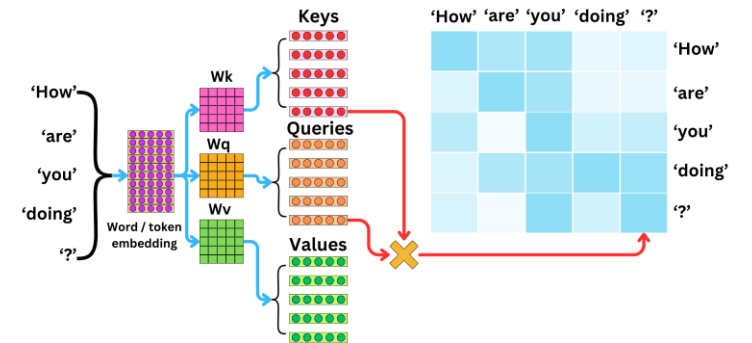
Se las hizo
interactuar entre
ellas



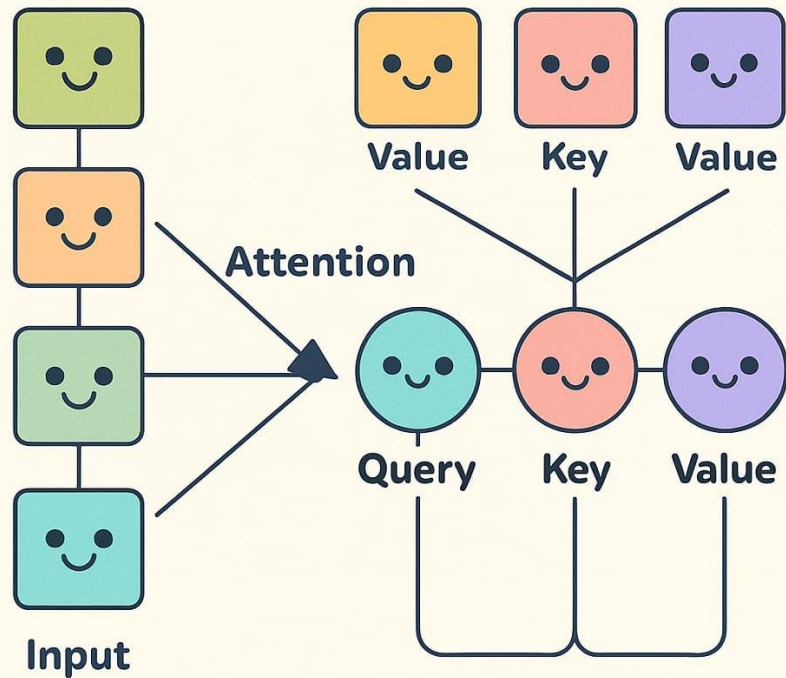
THE ATTENTION PLAYGROUND: TRANSFORMER NEURONAL NETWORKS!



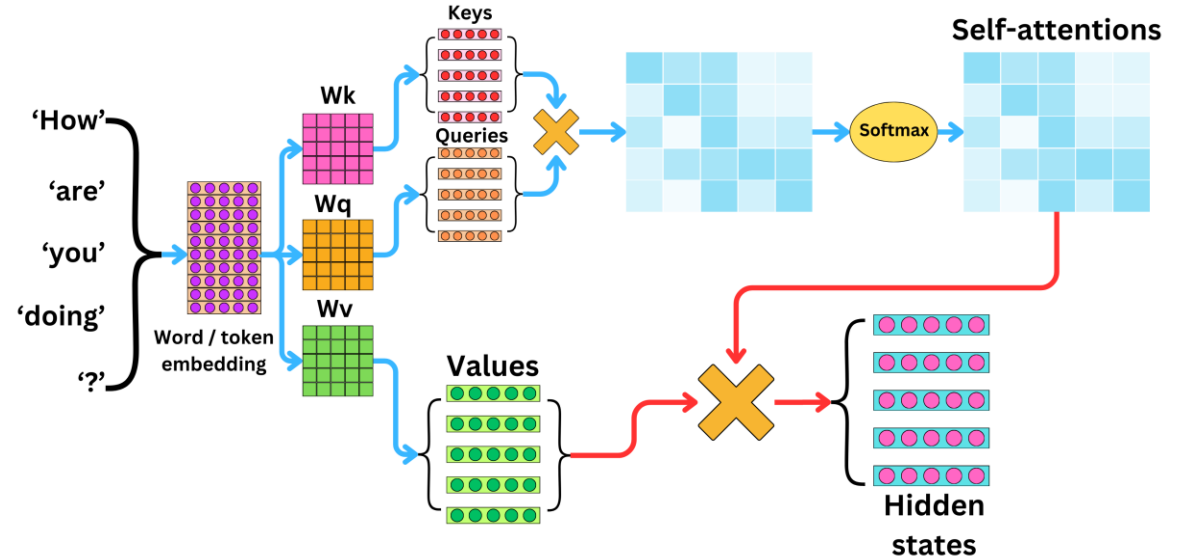
Mirarse



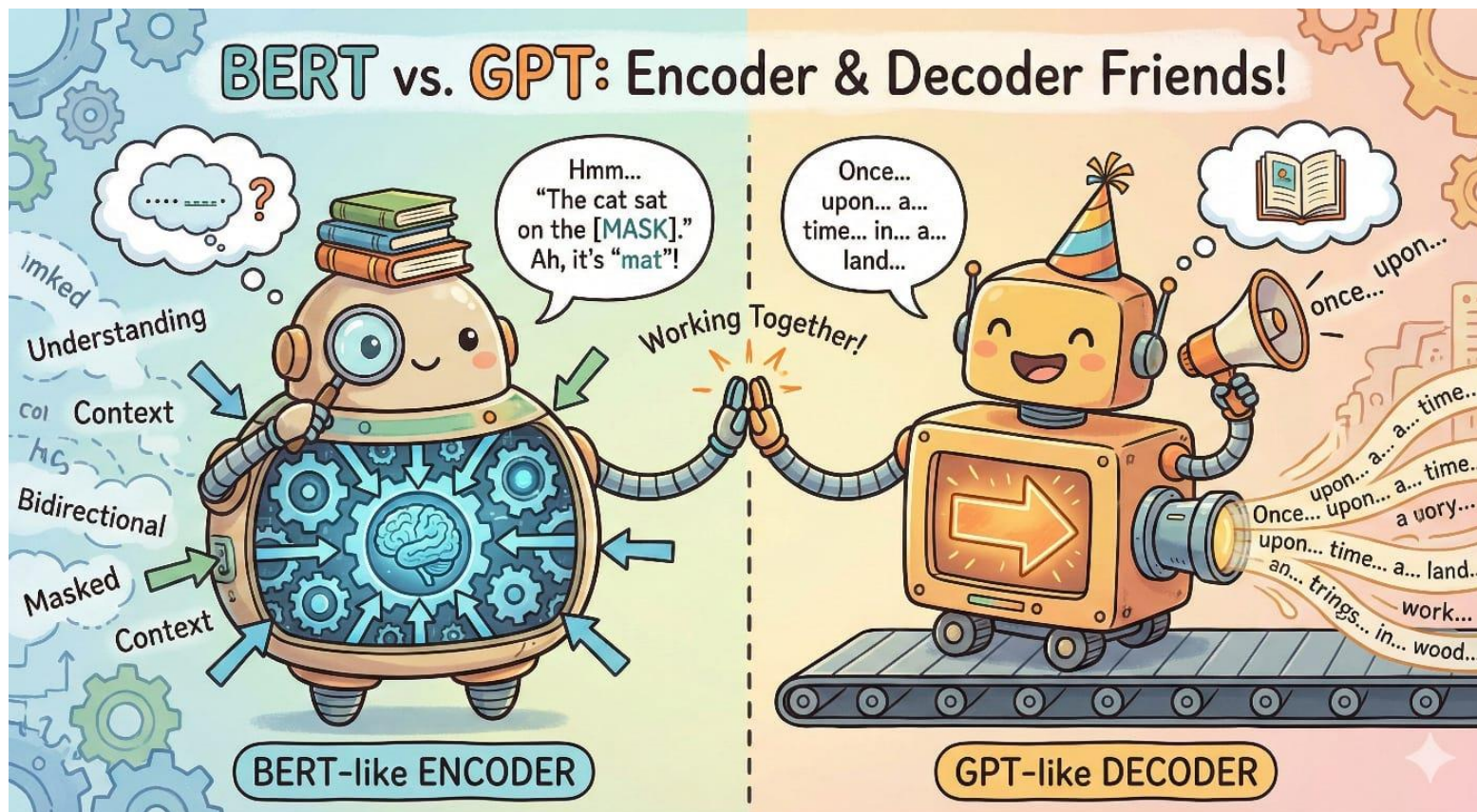
Transformer Attention



A través de complejos
mecanismos,
compararse a sí
mismas

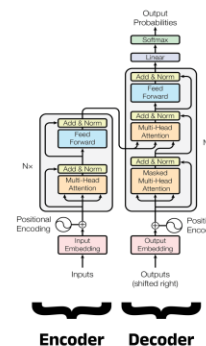


BERT vs. GPT: Encoder & Decoder Friends!

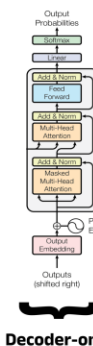


Se separó en
codificadores y
decodificadores de
información

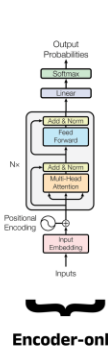
Transformer



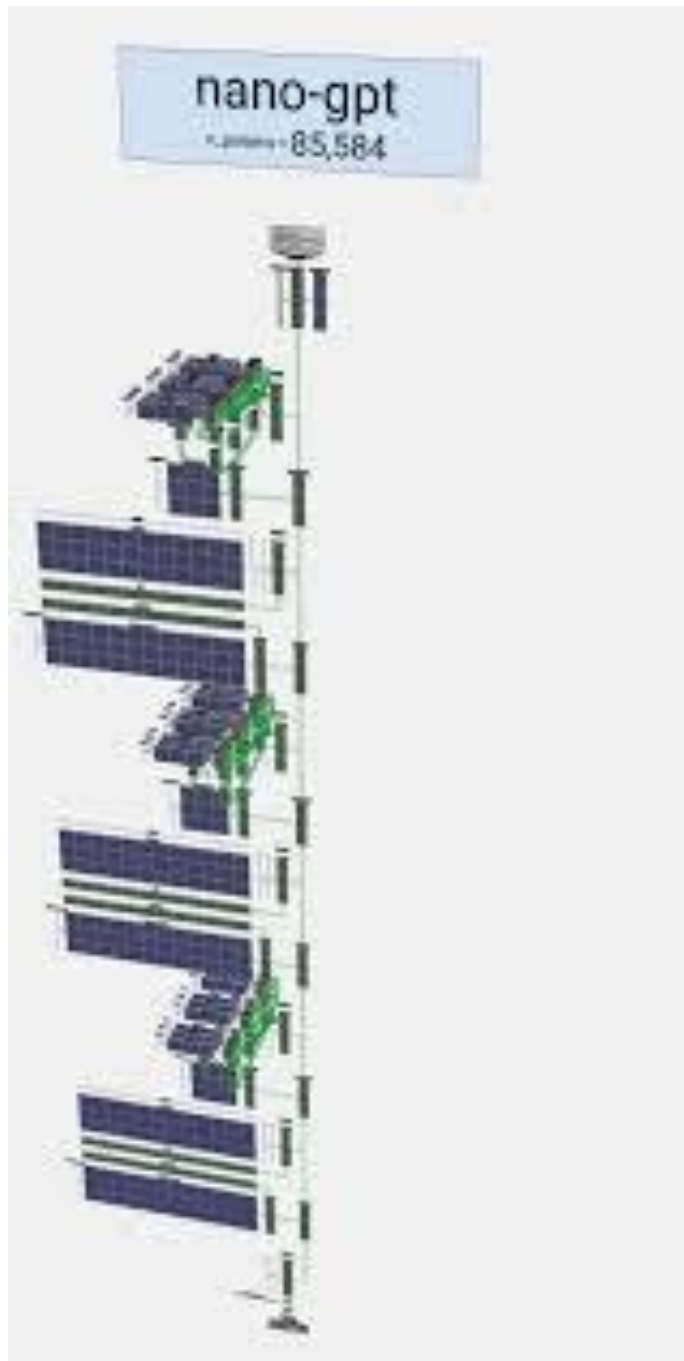
GPT*



BERT*

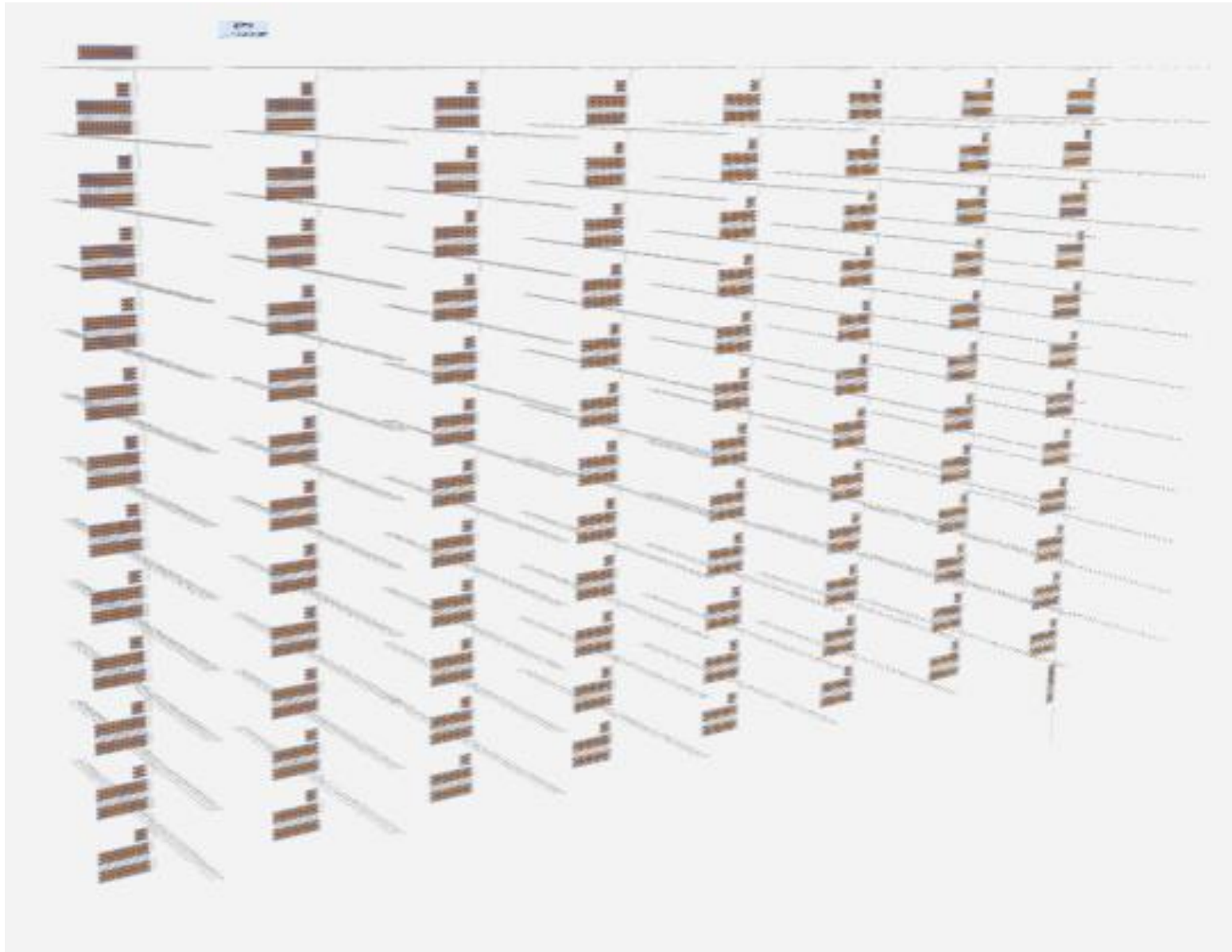


*Illustrative example, exact model architecture may vary slightly

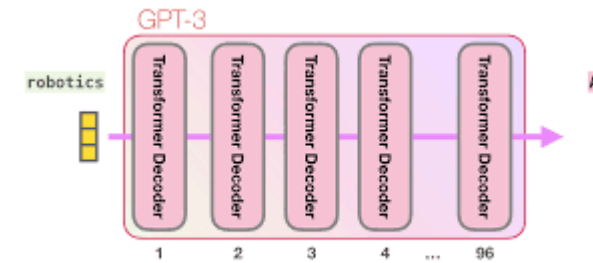


Se juntaron los
bloques y boom!
Nació un pequeño
GPT

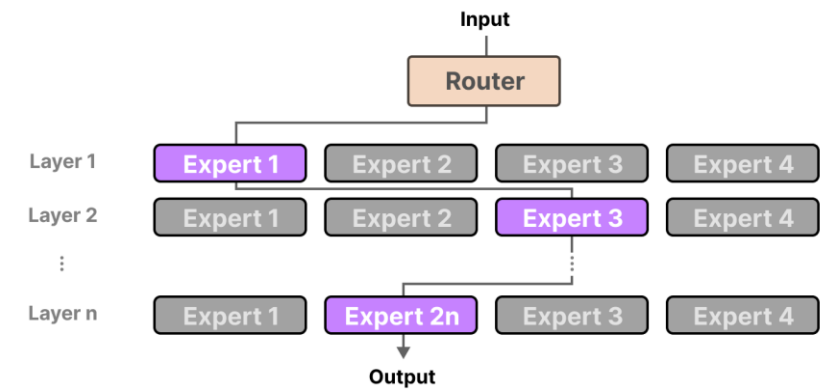
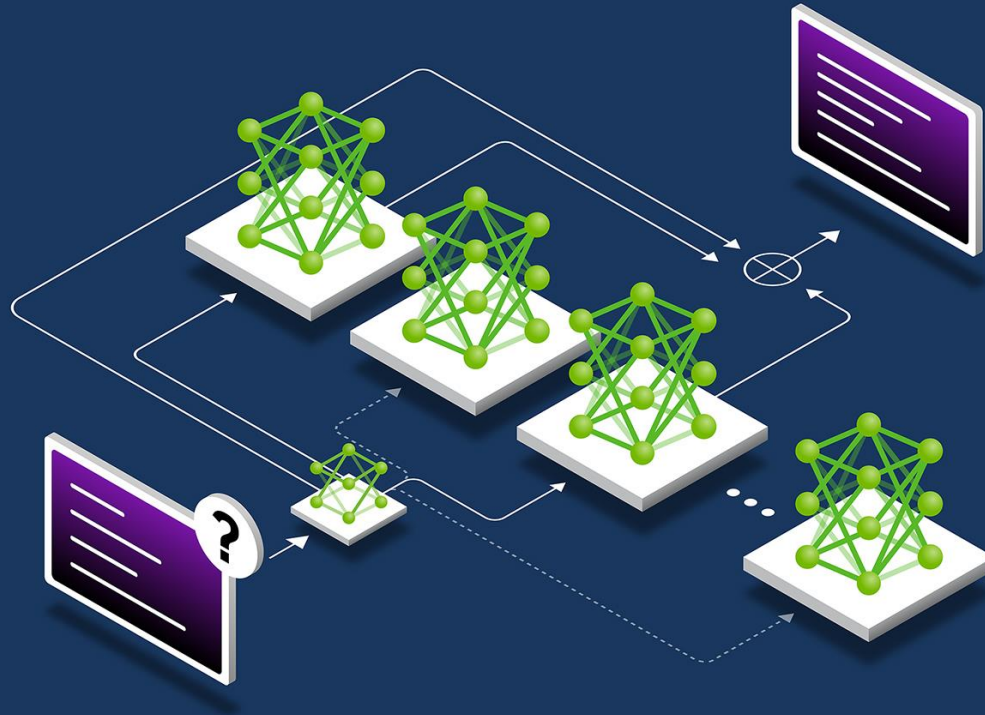
- NanoGPT tiene 3 bloques de 3 cabezales cada uno
- Son 85mil neuronas!



Y más y más
grandes... 96 bloques
de 96 cabezales cada
uno... y nació GPT3



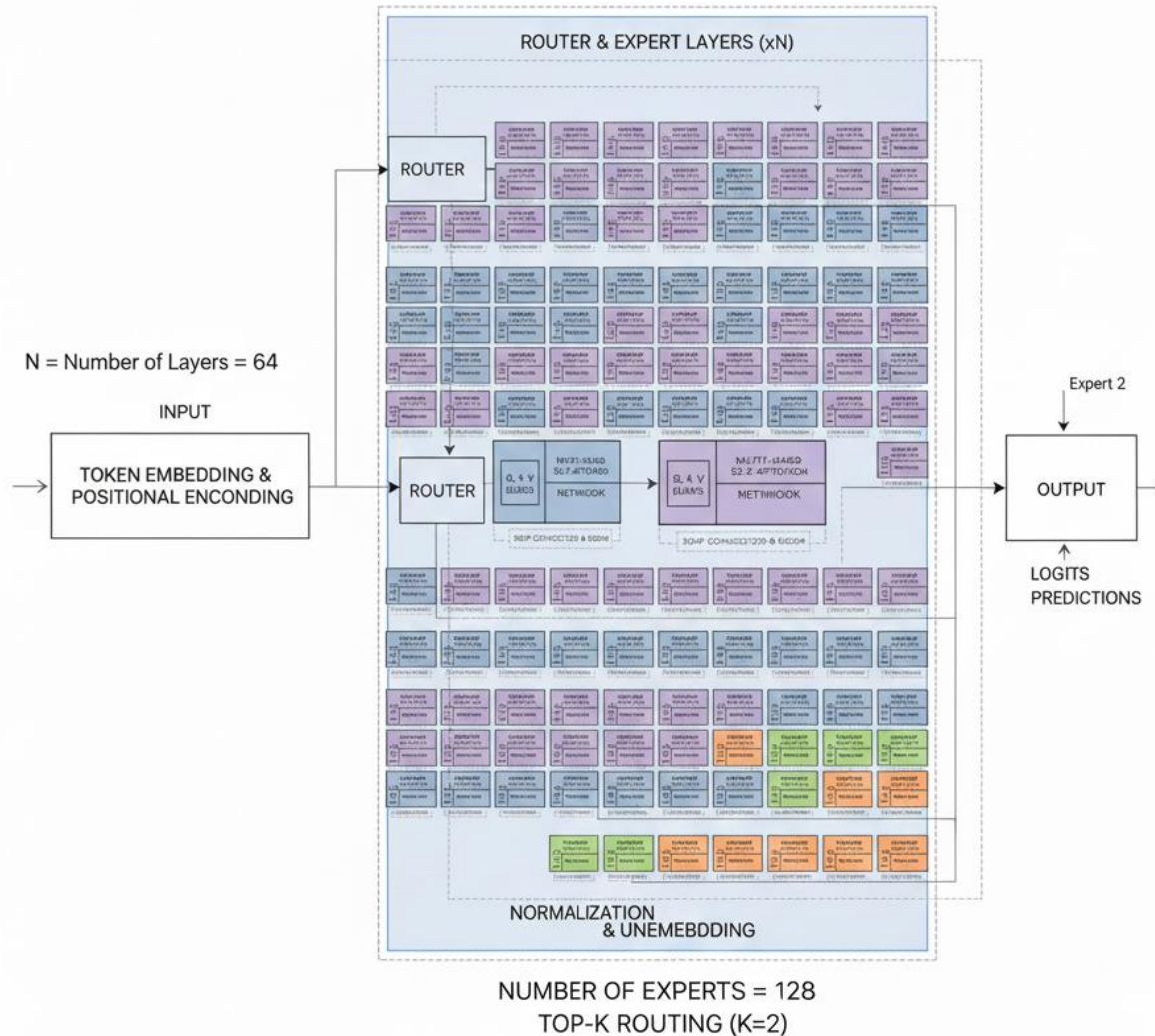
Ahora tienen cientos
y miles de bloques
con sub redes



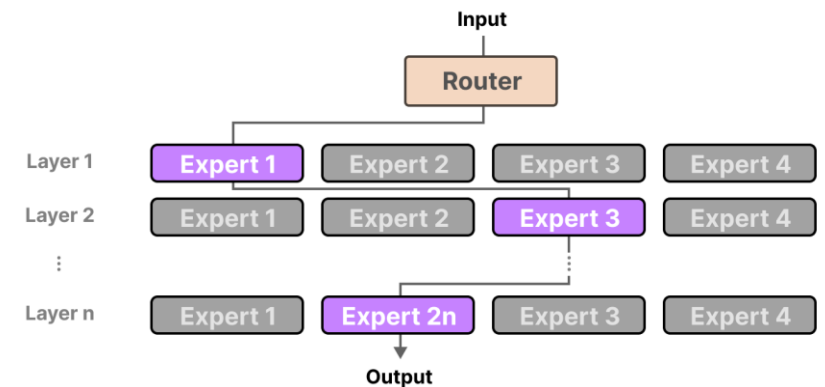
DEEPSEEK-MoE: GIANT EXPERT ARCHITECTURE

TOTAL PARAMETERS: 236 BILLION

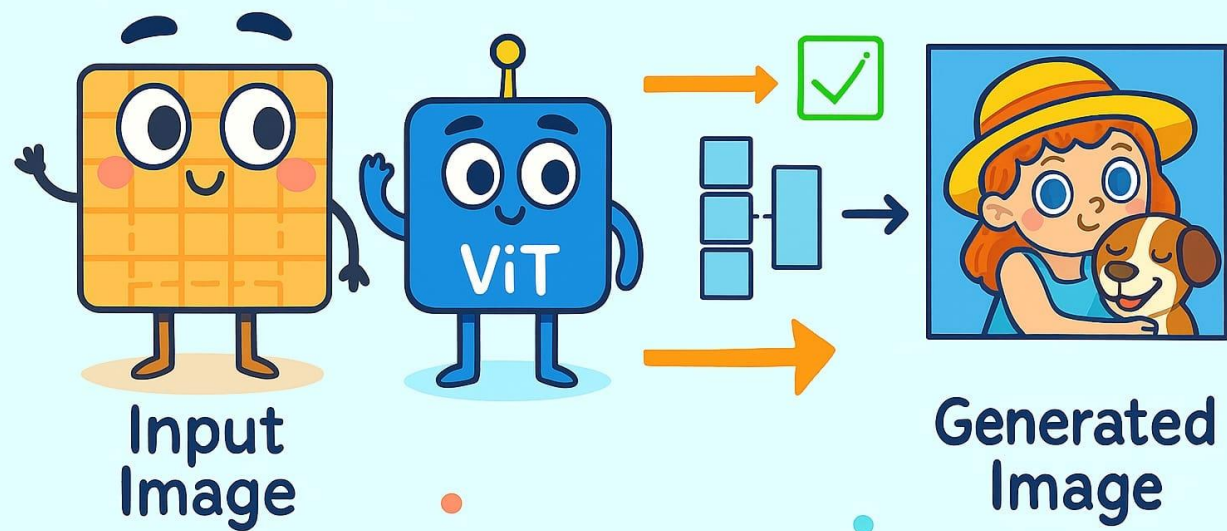
2 Trillion Tokens Pre-trained



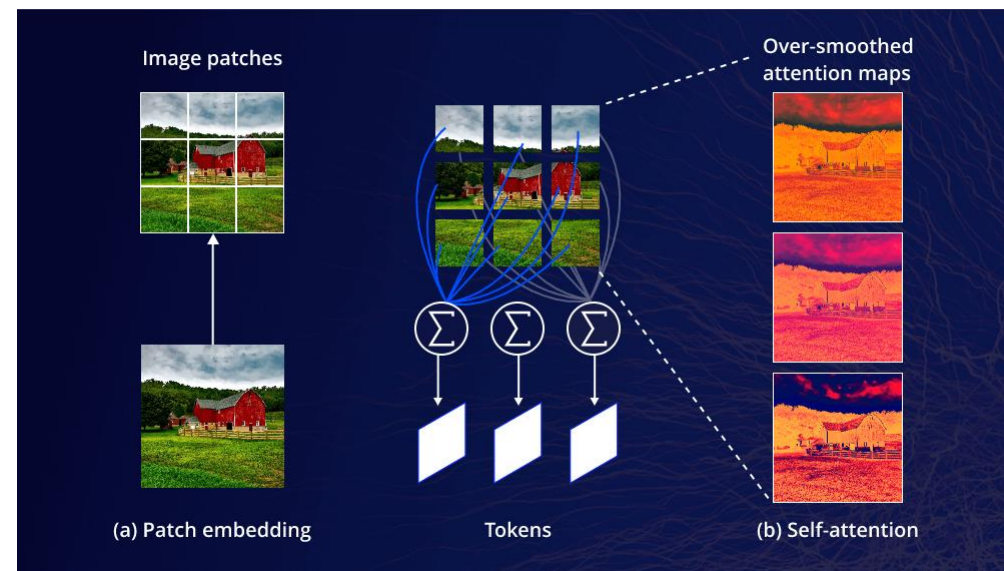
Cientos de expertos,
cada uno con decenas de
bloques con decenas de
cabezas cada bloque

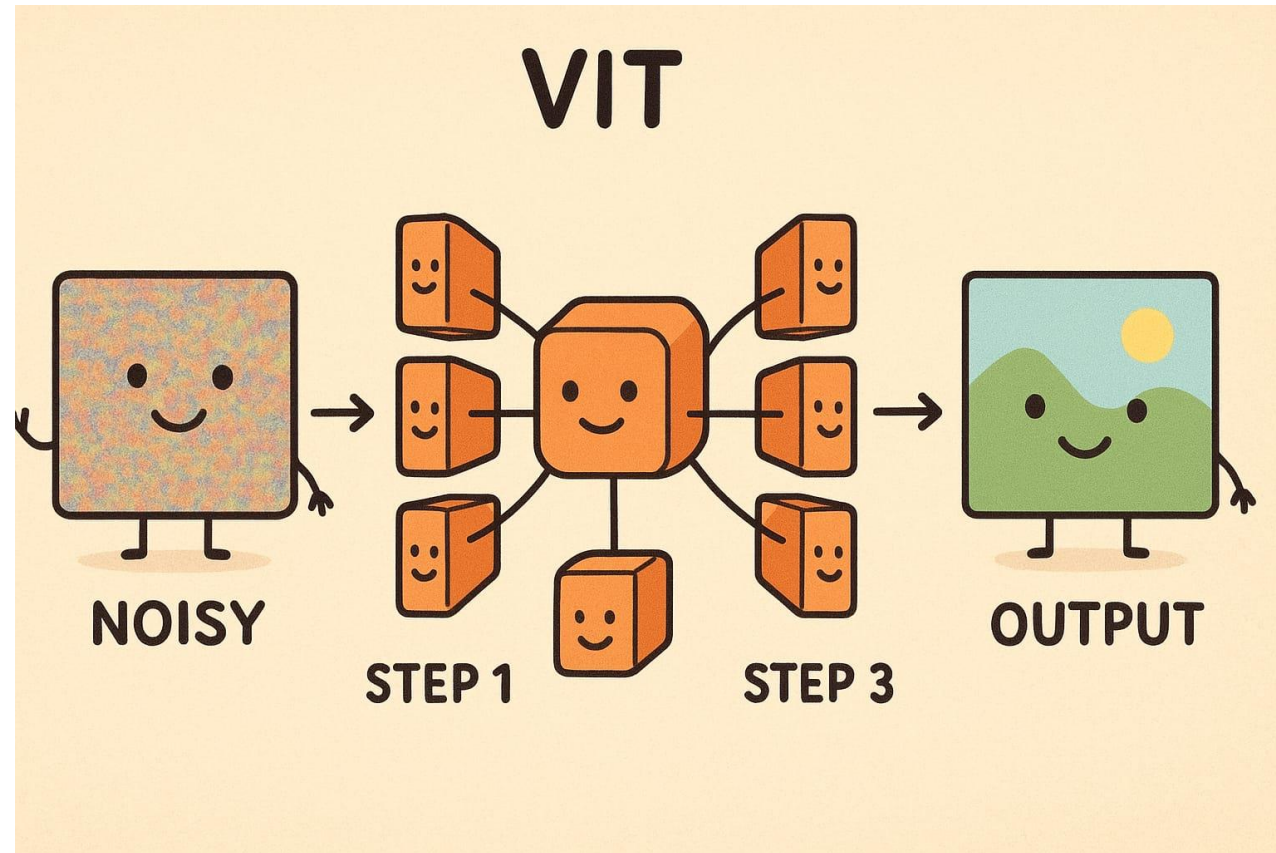


Vision Transformer Diffusion

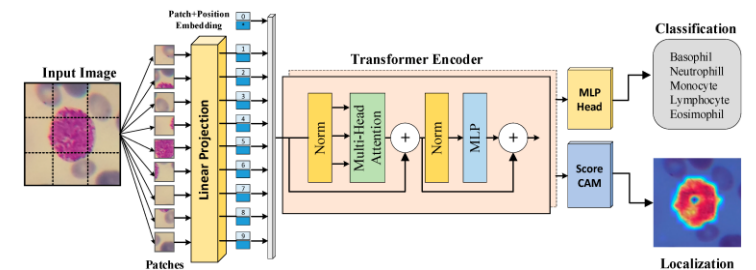


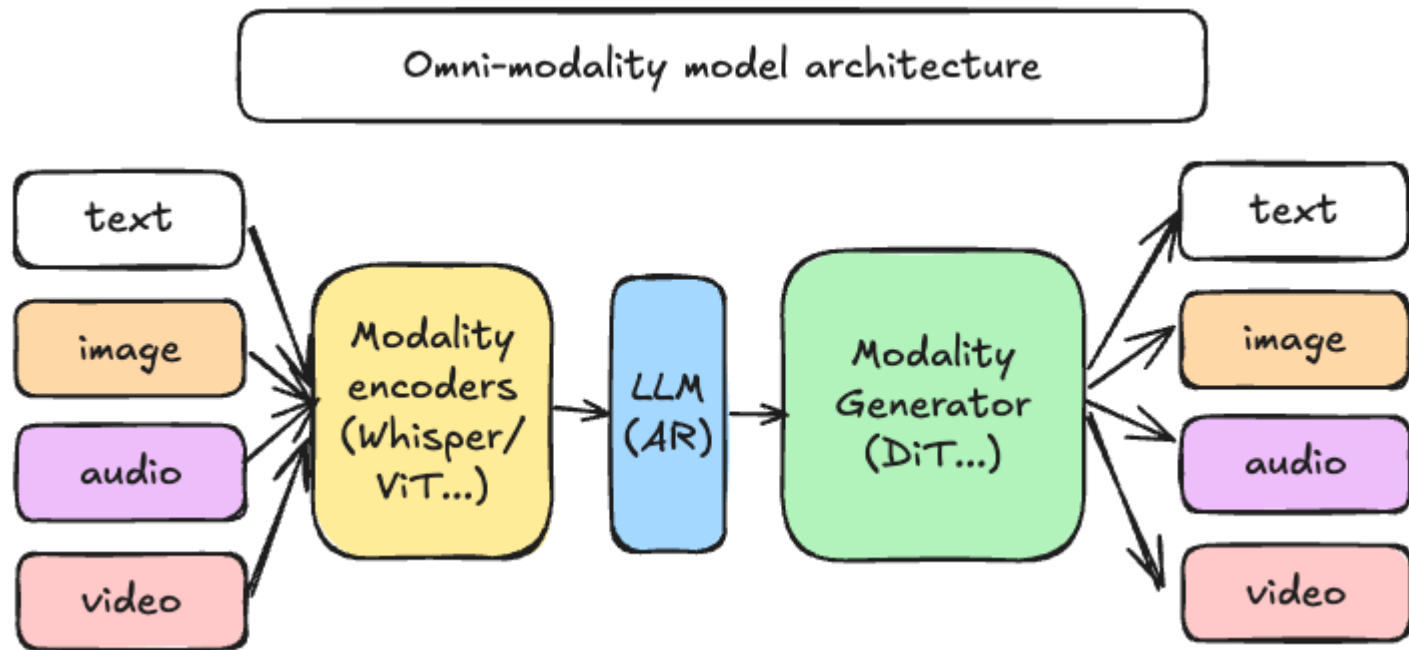
Y de pronto no sólo
era texto, también
imágenes



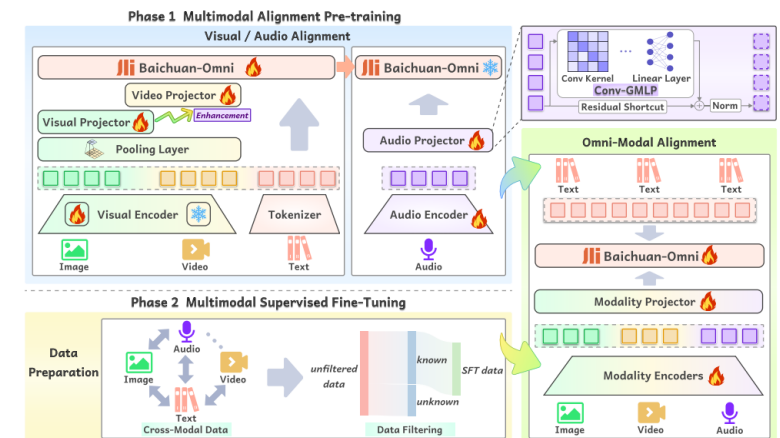


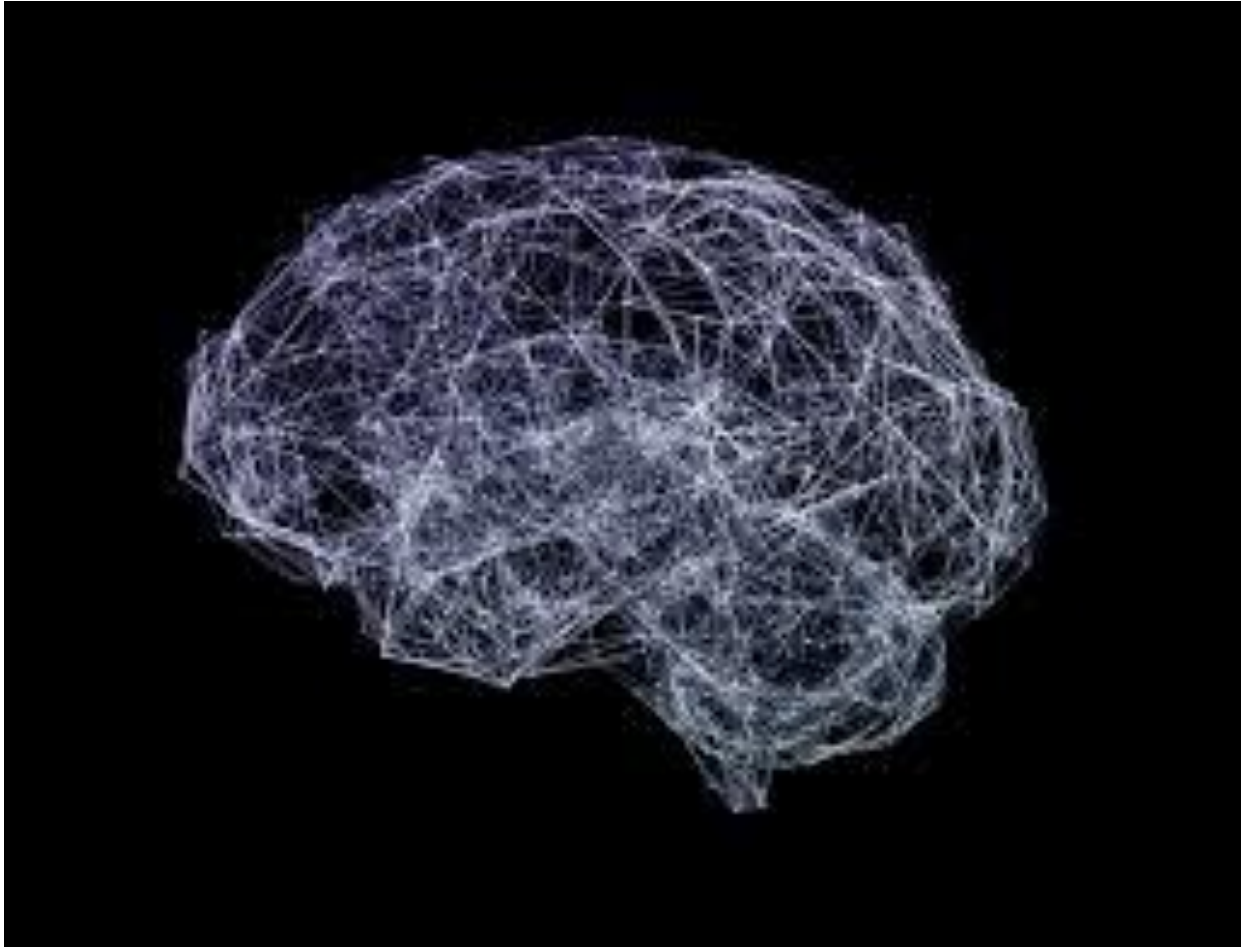
Y de pronto fuimos capaces de generar imágenes más y más perfectas





Finalmente, tenemos redes neuronales capaces de entender y generar de todo





El cerebro tiene entre 100T a 1000T de conexiones, conectadas en paralelo.

Los mejores modelos llegan a 1T o 2T, quizás a 5T los más caros, con conexiones secuenciales en cadena.

Aún nos falta mucho!