# M.EIC

# Natural Language Processing

Henrique Lopes Cardoso
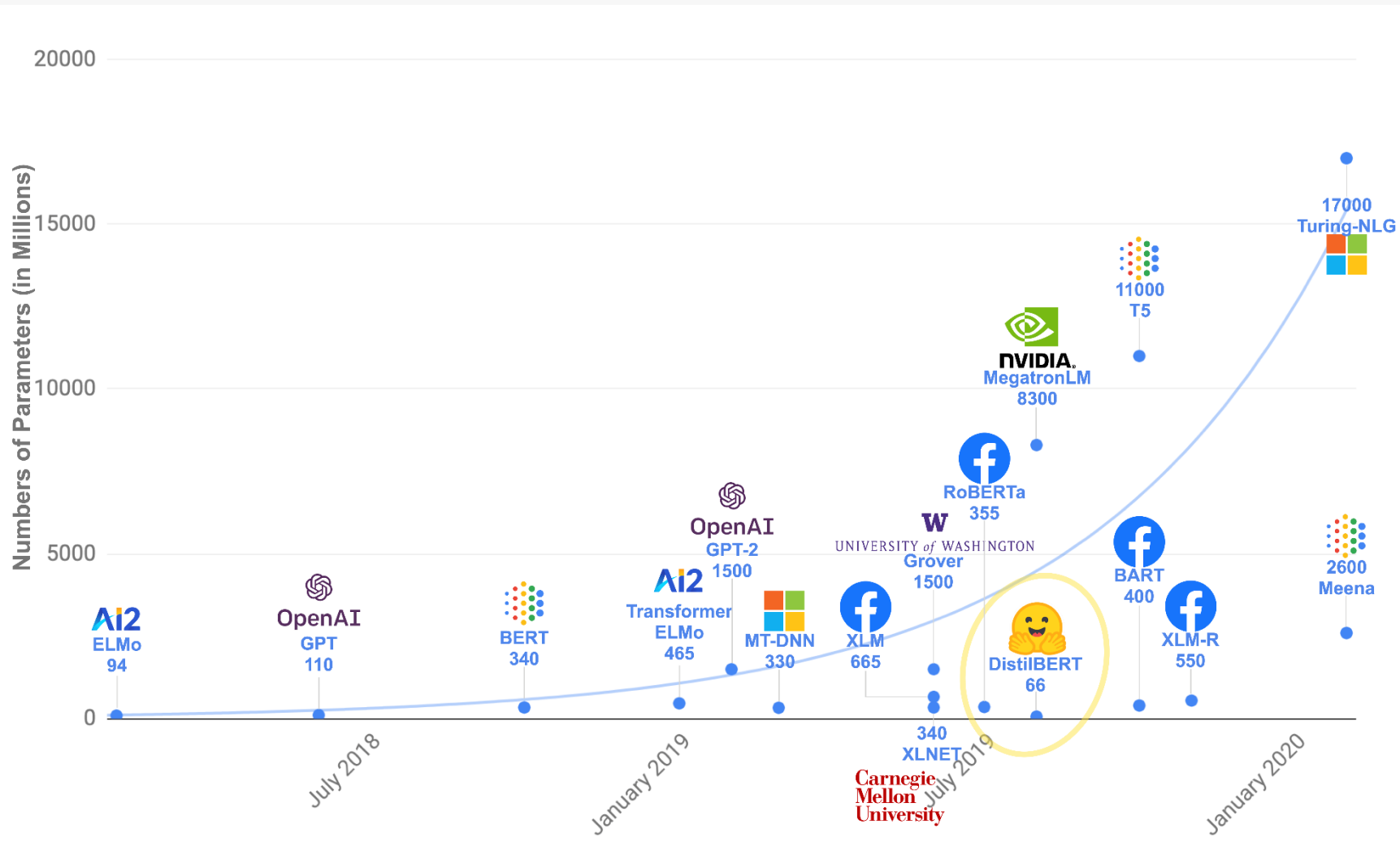
FEUP / LIACC

hlc@fe.up.pt

# Hugging Face Transformers



https://huggingface.co/
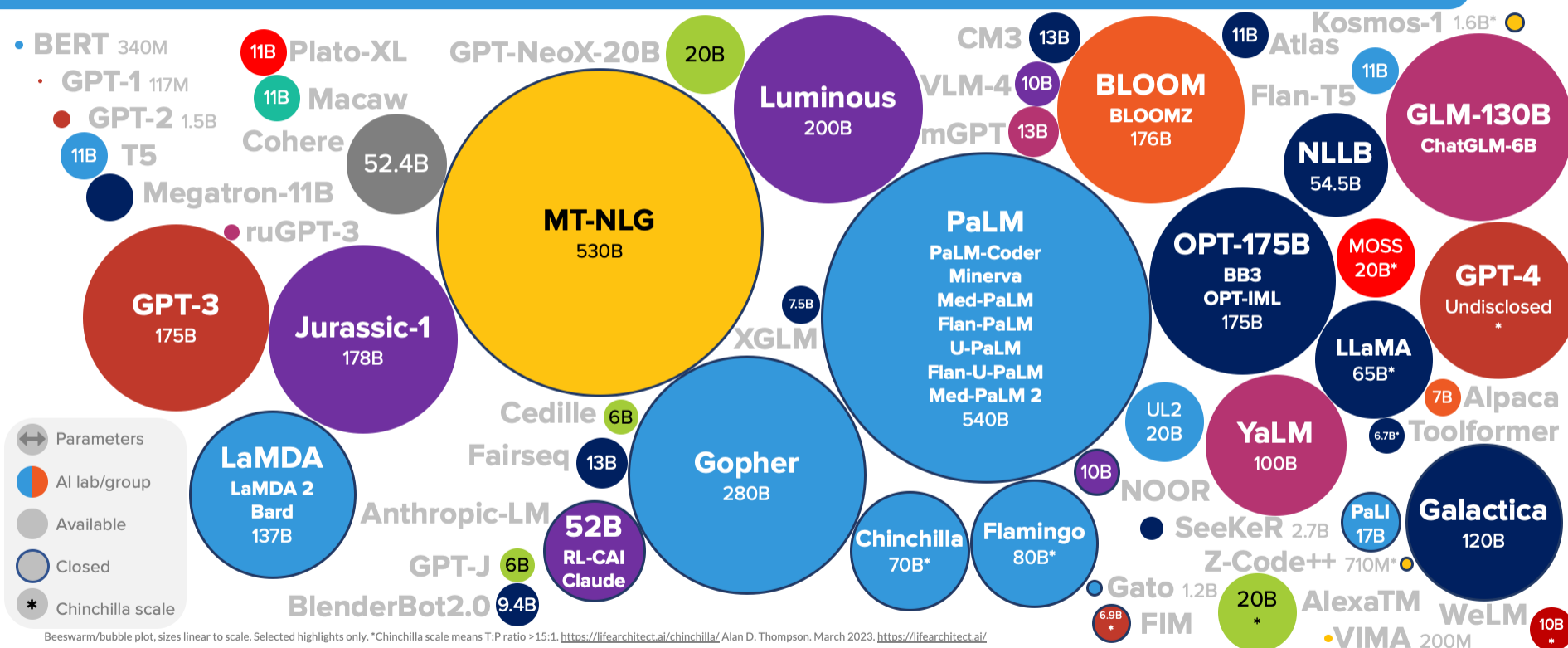
# Transformers are big models

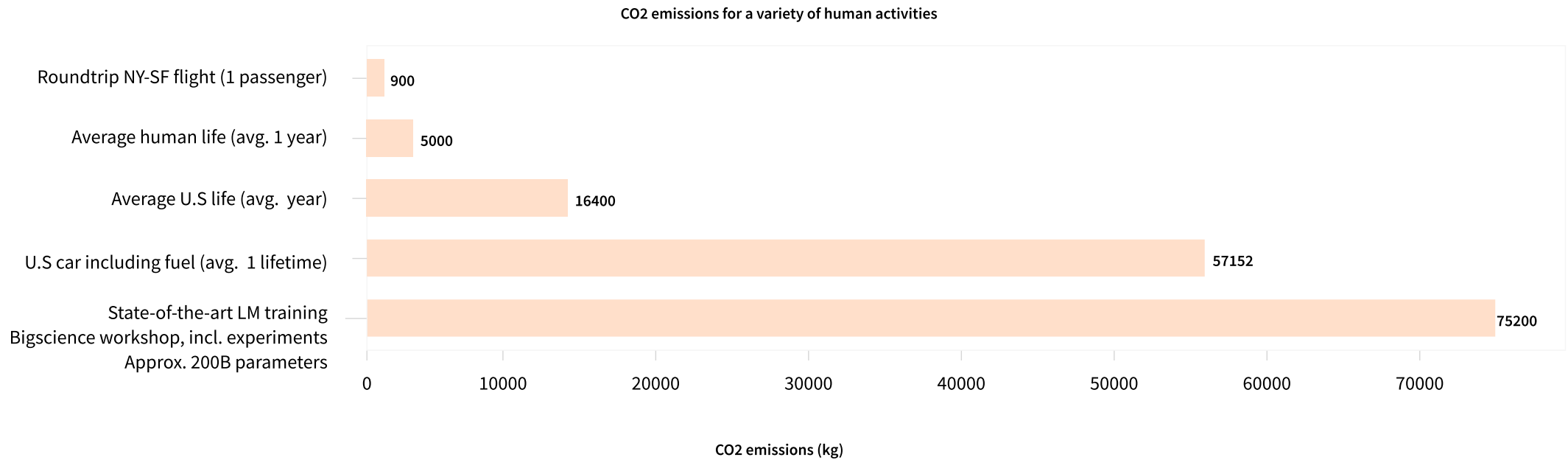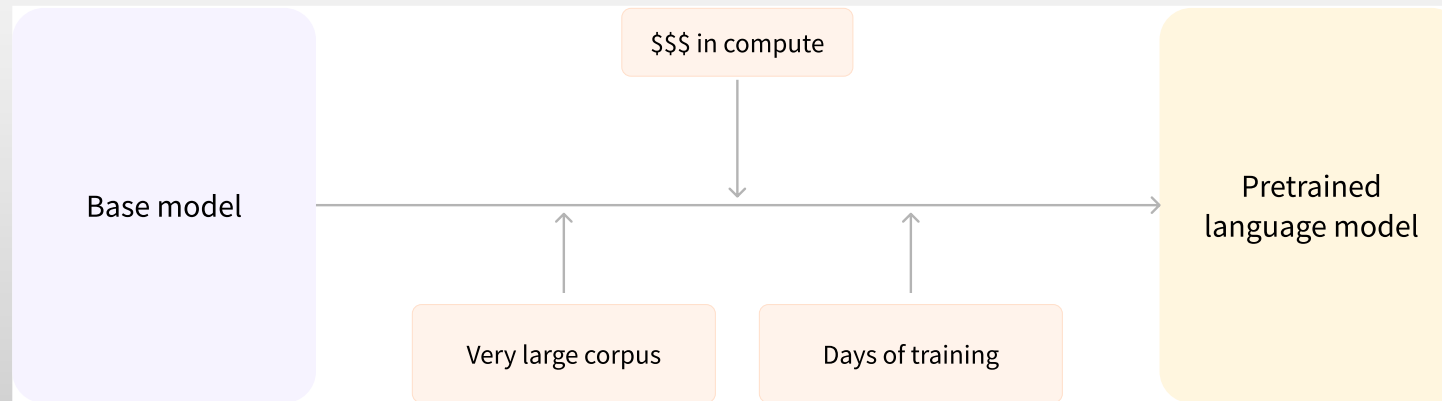# Transformers are big models



LANGUAGE MODEL SIZES TO MAR/2023

BERT 340M · GPT-1 117M · GPT-2 1.5B · Plato-XL 11B · Macaw 11B · Cohere · T5 11B · Megatron-11B · ruGPT-3 · GPT-NeoX-20B 20B · 52.4B · MT-NLG 530B · Luminous 200B · CM3 13B · VLM-4 10B · mGPT 13B · BLOOM BLOOMZ 176B · Atlas 11B · Kosmos-1 1.6B* · Flan-T5 11B · GLM-130B ChatGLM-6B · NLLB 54.5B · GPT-3 175B · Jurassic-1 178B · XGLM 7.5B · PaLM PaLM-Coder Minerva Med-PaLM Flan-PaLM U-PaLM Flan-U-PaLM Med-PaLM 2 540B · OPT-175B BB3 OPT-IML 175B · MOSS 20B* · GPT-4 Undisclosed * · LLaMA 65B* · UL2 20B · YaLM 100B · 6.7B* · Alpaca 7B · Toolformer · Cedille 6B · Fairseq 13B · Gopher 280B · 10B · NOOR · SeeKeR 2.7B · PaLI 17B · Galactica 120B · LaMDA LaMDA 2 Bard 137B · Anthropic-LM · 52B RL-CAI Claude · Chinchilla 70B* · Flamingo 80B* · Z-Code++ 710M* · GPT-J 6B · BlenderBot2.0 9.4B · Gato 1.2B · 6.9B* · FIM · 20B * · AlexaTM · WeLM · 10B · VIMA 200M

Parameters · AI lab/group · Available · Closed · Chinchilla scale

Beeswarm/bubble plot, sizes linear to scale. Selected highlights only. *Chinchilla scale means T:P ratio >15:1. https://lifearchitect.ai/chinchilla/ Alan D. Thompson. March 2023. https://lifearchitect.ai/

🔗 LifeArchitect.ai/models

# The carbon footprint

**CO2 emissions for a variety of human activities**

| Activity | CO2 emissions (kg) |
|---|---|
| Roundtrip NY-SF flight (1 passenger) | 900 |
| Average human life (avg. 1 year) | 5000 |
| Average U.S life (avg. year) | 16400 |
| U.S car including fuel (avg. 1 lifetime) | 57152 |
| State-of-the-art LM training Bigscience workshop, incl. experiments Approx. 200B parameters | 75200 |

CO2 emissions (kg)

# Transfer Learning

- Pretraining



- Fine-tuning

# Transformer Architecture

# Encoder, Decoder, Encoder-Decoder

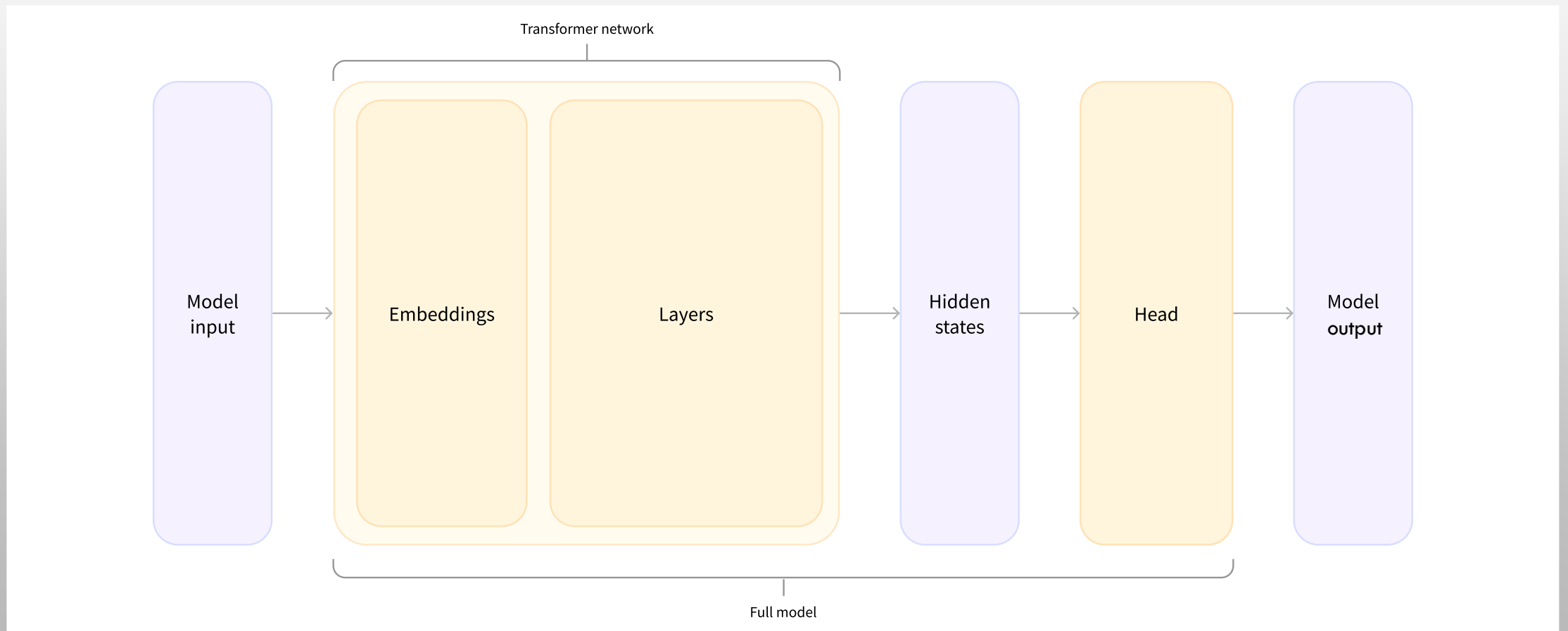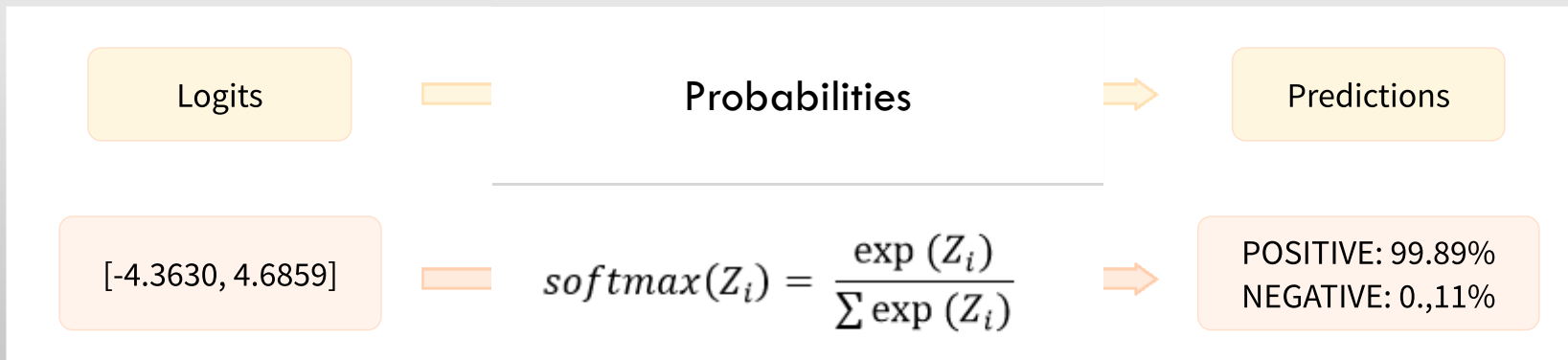| Model | Examples | Tasks |
|---|---|---|
| Encoder | ALBERT, BERT, DistilBERT, ELECTRA, RoBERTa | Sentence classification, named entity recognition, extractive question answering |
| Decoder | CTRL, GPT, GPT-2, Transformer XL | Text generation |
| Encoder-decoder | BART, T5, Marian, mBART | Summarization, translation, generative question answering |

# Pipelines

# Tokenizer

# Model

# Post Processing

Logits — Probabilities ⇒ Predictions

$[-4.3630, 4.6859]$ — $softmax(Z_i) = \dfrac{\exp(Z_i)}{\sum \exp(Z_i)}$ ⇒ POSITIVE: 99.89%
NEGATIVE: 0.,11%
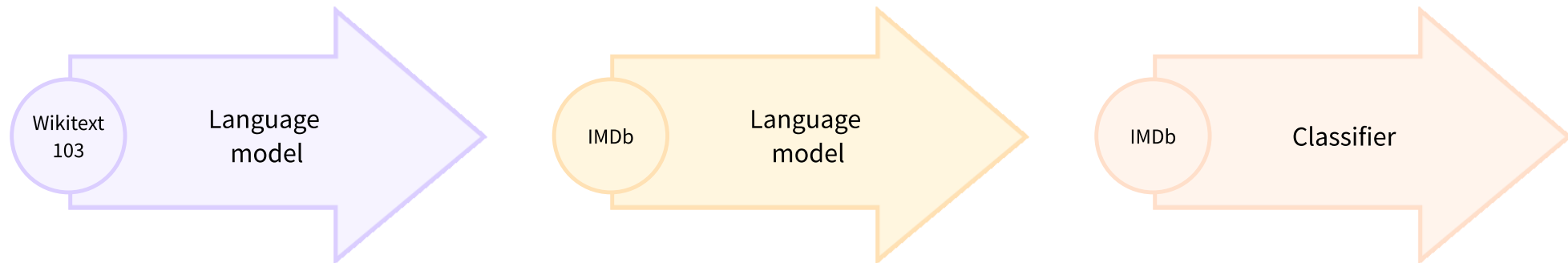
# NLP Tasks

- Sequence classification

- Token classification (sequence labelling)

- Masked language modeling (like BERT)

- Summarization

- Translation

- Causal language modeling pretraining (like GPT-2)

- Question answering

# Domain Adaptation (through MLM)

# Models

Tasks

| | | | |
|---|---|---|---|
| Image Classification | | Translation | |
| Unconditional Image Generation | | Fill-Mask | |
| Automatic Speech Recognition | | Token Classification | |
| Sentence Similarity | | Audio Classification | |
| Question Answering | | Summarization | |
| Zero-Shot Classification | + 15 | | |

Libraries

PyTorch    TensorFlow    JAX    + 24

Datasets

common_voice    wikipedia    squad    glue

bookcorpus    c4    conll2003    emotion    + 938

Languages

en    es    fr    de    zh    sv    ru    fi    + 173

Models  40,844          Search Models

**gpt2**
Text Generation · Updated May 19, 2021 · ↓ 59.9M · ♡ 83

**bert-base-uncased**
Fill-Mask · Updated May 18, 2021 · ↓ 16.6M · ♡ 136

**cross-encoder/ms-marco-MiniLM-L-12-v2**
Text Classification · Updated Aug 5, 2021 · ↓ 9.87M · ♡ 6

**distilbert-base-uncased-finetuned-sst-2-english**
Text Classification · Updated Mar 22 · ↓ 4.66M · ♡ 50

**Helsinki-NLP/opus-mt-zh-en**
Translation · Updated Feb 26, 2021 · ↓ 3.81M · ♡ 25

**sentence-transformers/all-MiniLM-L6-v2**
Sentence Similarity · Updated Aug 30, 2021 · ↓ 2.97M · ♡ 34

**distilgpt2**
Text Generation · Updated May 21, 202

**roberta-base**
Fill-Mask · Updated Jul 6, 2021 · ↓ 12

**distilbert-base-uncased**
Fill-Mask · Updated Aug 29, 2021 · ↓ 4

**xlm-roberta-large-finetune**
Token Classification · Updated Oct 12,

**bert-base-chinese**
Fill-Mask · Updated May 18, 2021 · ↓

**bert-base-cased**
Fill-Mask · Updated Sep 6, 2021 · ↓ 2.

# Datasets

# Explore 🤗

- Hugging Face documentation: https://huggingface.co/docs

- Hugging Face Transformers documentation: https://huggingface.co/docs/transformers/

- Hugging Face course: https://huggingface.co/course/

- Hugging Face models: https://huggingface.co/models

- Hugging Face datasets: https://huggingface.co/datasets