

Transformers

The essentials of Hugging Face Transformers, from one-call pipelines to tokenization, text generation, and fine-tuning with Trainer.

Install & Load from the Hub

Install transformers, then pull any of 1M+ models by name.

PURPOSE	COMMAND	RESULT
Install transformers + torch.	<code>pip install "transformers[torch]"</code>	
Check the installed version.	<code>import transformers; transformers.__version__</code>	
Log in to the Hub (private/push).	<code>hf auth login</code>	
Load a model by Hub id.	<code>AutoModel.from_pretrained("bert-base-uncased")</code>	
Load weights in their stored dtype.	<code>AutoModel.from_pretrained(name, dtype="auto")</code>	
Spread a big model across devices.	<code>AutoModelForCausalLM.from_pretrained(name, device_map="auto")</code>	

`from_pretrained("org/model")` downloads and caches once; `dtype="auto"` loads stored precision, `device_map="auto"` shards big models

pipeline(): a task in 3 lines

One call: name a task, get a ready-to-run model.

PURPOSE	COMMAND	RESULT
Build a task pipeline.	<code>pipe = pipeline(task="sentiment-analysis")</code>	
Run it on text.	<code>pipe("I love this!")</code>	
Pick a specific model.	<code>pipeline("text-generation", model="Qwen/Qwen3-0.6B")</code>	
Choose the device.	<code>pipeline(task, model=name, device="cuda")</code>	
Batch a list of inputs.	<code>pipe(["text a", "text b", "text c"], batch_size=8)</code>	
Generation kwargs pass through.	<code>pipe(prompt, max_new_tokens=50, do_sample=True)</code>	

`pipeline(task=...)` loads a default model; set `model=` and `device=` to choose, pass a list with `batch_size=`, and extra kwargs flow to the model

Tokenize

The tokenizer turns text into integer ids the model reads.

PURPOSE	COMMAND	RESULT
Load the matching tokenizer.	<code>tok = AutoTokenizer.from_pretrained(name)</code>	
Encode to tensors.	<code>enc = tok("Hello world", return_tensors="pt")</code>	
See the integer ids.	<code>enc["input_ids"]</code>	
Read the attention mask.	<code>enc["attention_mask"]</code>	
Decode ids back to text.	<code>tok.decode(enc["input_ids"][0])</code>	
Apply a chat template.	<code>tok.apply_chat_template(messages, add_generation_prompt=True, return_tensors="pt")</code>	

load AutoTokenizer from the same checkpoint as the model; tok(text, return_tensors="pt") gives input_ids + attention_mask, tok.decode goes back

AutoModel & AutoTokenizer







Auto-classes infer the right architecture from the name.

PURPOSE	COMMAND	RESULT
Raw backbone (hidden states only).	<code>AutoModel.from_pretrained(name)</code>	
Causal LM head (text generation).	<code>AutoModelForCausalLM.from_pretrained(name)</code>	
Sequence classification head.	<code>AutoModelForSequenceClassification.from_pretrained(name, num_labels=2)</code>	
Token classification (NER).	<code>AutoModelForTokenClassification.from_pretrained(name, num_labels=9)</code>	
Pair tokenizer to model.	<code>tok = AutoTokenizer.from_pretrained(name)</code>	
Save both locally.	<code>model.save_pretrained("out/")</code> <code>tok.save_pretrained("out/")</code>	

AutoModel is the bare backbone; AutoModelFor* snaps on a task head; always pair AutoTokenizer on the same checkpoint

Batch, Pad & Truncate


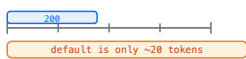

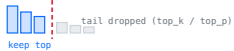

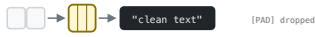
Encode many texts at once into one rectangular tensor.

PURPOSE	COMMAND	RESULT
Encode a list of texts.	<code>tok(["short", "a longer one"], padding=True)</code>	
Pad to a fixed length.	<code>tok(texts, padding="max_length", max_length=128)</code>	
Truncate long inputs.	<code>tok(texts, truncation=True, max_length=512)</code>	
Mask hides the padding.	<code>enc["attention_mask"]</code>	
Pad on the left for generation.	<code>AutoTokenizer.from_pretrained(name, padding_side="left")</code>	
Dynamic padding per batch.	<code>DataCollatorWithPadding(tokenizer=tok)</code>	

padding=True pads to the longest row; the attention_mask makes padding safe; pad left for decoder LLMs, collate dynamically to save compute

Generate Text







model.generate() decodes new tokens; dials control style.

PURPOSE	COMMAND	RESULT
Greedy: always pick the top token.	<code>model.generate(**enc, max_new_tokens=50)</code>	
Always set the output length.	<code>model.generate(**enc, max_new_tokens=200)</code>	
Sample for creative text.	<code>model.generate(**enc, do_sample=True, temperature=0.8)</code>	
Trim the candidate pool.	<code>model.generate(**enc, do_sample=True, top_k=50, top_p=0.95)</code>	
Curb repetition.	<code>model.generate(**enc, repetition_penalty=1.2)</code>	
Decode the new tokens to text.	<code>tok.batch_decode(out, skip_special_tokens=True)</code>	

generate() is greedy by default; always set max_new_tokens, then do_sample + temperature/top_k/top_p shape the output

Forward Pass & Logits

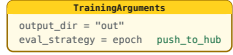

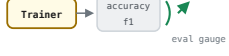

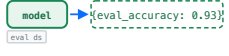

Call the model directly to read raw logits and hidden states.

PURPOSE	COMMAND	RESULT
Run a forward pass (no grad).	<code>with torch.no_grad(): out = model(**enc)</code>	
Read the logits.	<code>out.logits</code>	
Turn logits into probabilities.	<code>out.logits.softmax(dim=-1)</code>	
Take the predicted class.	<code>out.logits.argmax(dim=-1)</code>	
Ask for hidden states.	<code>model(**enc, output_hidden_states=True).hidden_states</code>	
Pooled sentence embedding.	<code>out.last_hidden_state[:, 0]</code>	

`model(**enc)` returns `.logits` (pre-softmax); `softmax` for probabilities, `argmax` for the class, `output_hidden_states=True` for per-layer states

Fine-tune with Trainer & push

Trainer runs the loop; `push_to_hub()` shares the result.

PURPOSE	COMMAND	RESULT
Configure the run.	<code>args = TrainingArguments(output_dir="out", eval_strategy="epoch", push_to_hub=True)</code>	
Assemble the trainer.	<code>trainer = Trainer(model, args, train_dataset, eval_dataset, processing_class=tok)</code>	
Score during training.	<code>Trainer(..., compute_metrics=compute_metrics)</code>	
Train the model.	<code>trainer.train()</code>	
Evaluate on the held-out set.	<code>trainer.evaluate()</code>	
Push weights + tokenizer to Hub.	<code>trainer.push_to_hub()</code>	

Trainer wraps the training loop: configure with `TrainingArguments`, pass the tokenizer as `processing_class=`, then `train()`, `evaluate()`, `push_to_hub()`