

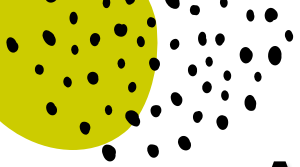
Intelligent User Interface

By Dr. Taghinezhad

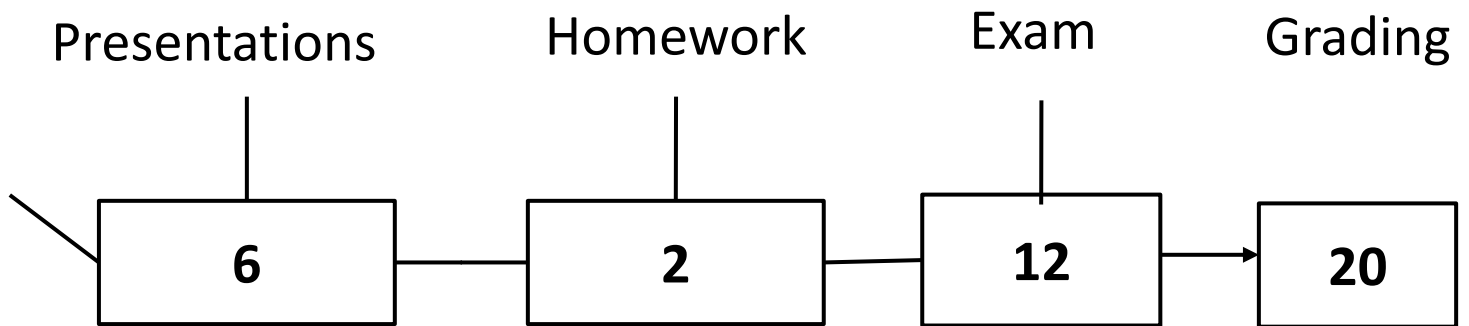
Mail:

a0taghinezhad@gmail.com

<https://ataghinezhad.github.io/>



About This Course



<https://ataghinezhad.github.io/>

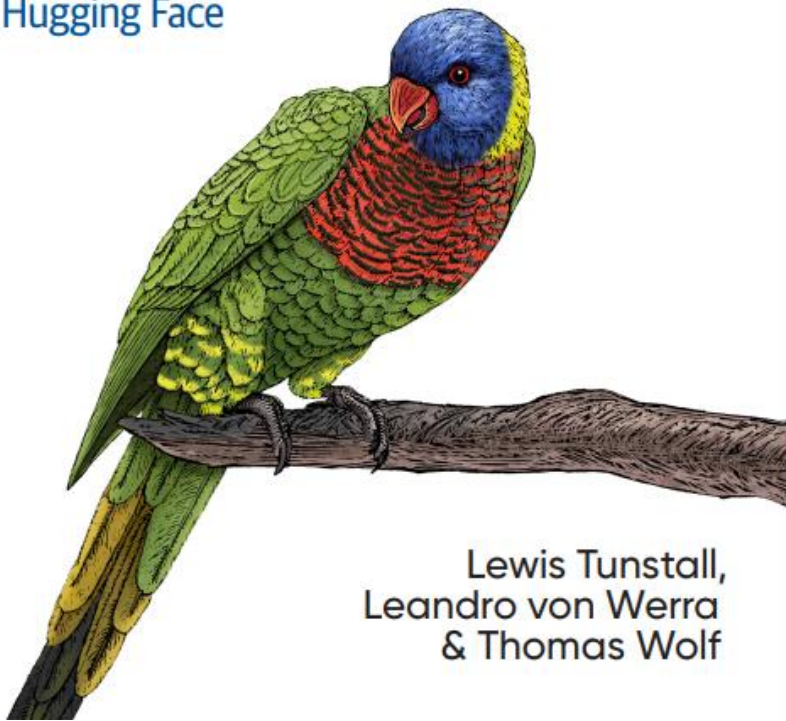
References

O'REILLY®

Revised
Edition

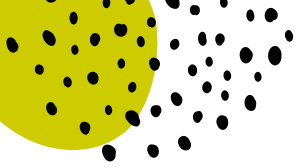
Natural Language Processing with Transformers

Building Language Applications
with Hugging Face



Lewis Tunstall,
Leandro von Werra
& Thomas Wolf

d.github.io/



Projects include:

Each student must review at least three papers and write a literature review on these subjects. The subjects are as follows, but not limited to them:

1. Integration of Reinforcement Learning with Transformer Architectures for Adaptive Interfaces:

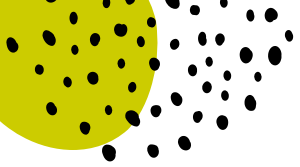
1. Real-Time User Behavior Modeling:
2. Dynamic Interface Optimization:

2. Delve into the ethical implications of deploying transformer models in user interfaces:

1. Bias Mitigation and Transparency: Evaluate how biases in transformer outputs can be identified and mitigated in real-world applications.
2. User Privacy and Data Security: Study the trade-offs between personalization and privacy, proposing guidelines for responsible interface de

3. Conversational Agents and Dialogue Systems Examine how transformer-based models can be deployed in the development of conversational user interfaces:

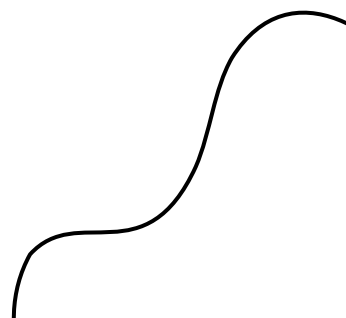
1. Contextual Dialogue Management: Research advanced techniques for maintaining coherent, context-aware conversations over extended interactions.
2. Emotional and Sentiment Analysis:

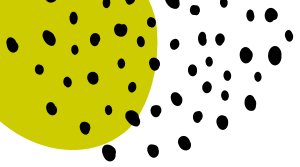


Chapters:

- CHAPTER 1 Hello Transformers
- CHAPTER 2 Text Classification
- CHAPTER 3 Transformer Anatomy
- CHAPTER 5 Text Generation
- CHAPTER 7 Question Answering

<https://ataghinezhad.github.io/>





Chapters:

- Chapter 1, Hello Transformers, introduces transformers and puts them into context. It also provides an introduction to the Hugging Face ecosystem.
- Chapter 2, Text Classification, focuses on the task of sentiment analysis (a common text classification problem) and introduces the Trainer API.
- Chapter 3, Transformer Anatomy, dives into the Transformer architecture in more depth, to prepare you for the chapters that follow.
- Chapter 5, Text Generation, explores the ability of transformer models to generate text, and introduces decoding strategies and metrics.
- Chapter 7, Question Answering, focuses on building a review-based question answering system and introduces retrieval with Haystack.

CHAPTER 1 Hello Transformers

The Rise of Transformer Models in NLP

- **Introduction**
- **Key Advances in NLP (2017):**
 - **Transformer Architecture** (Google): Superior to RNNs for translation tasks.
 - **ULMFiT**: Transfer learning with LSTMs for state-of-the-art text classification.
- **Impact:** Foundations for GPT and BERT models.

<https://ataghinezhad.github.io/>

Introduction

- **Key Advances in NLP (2017):**

- Transformer Architecture (Google)
- ULMFiT for text classification

- **Impact:** Laid the groundwork for GPT, BERT, and a wide range of transformer models

- **Reference:** Timeline of transformer models (see Figure 1-1)

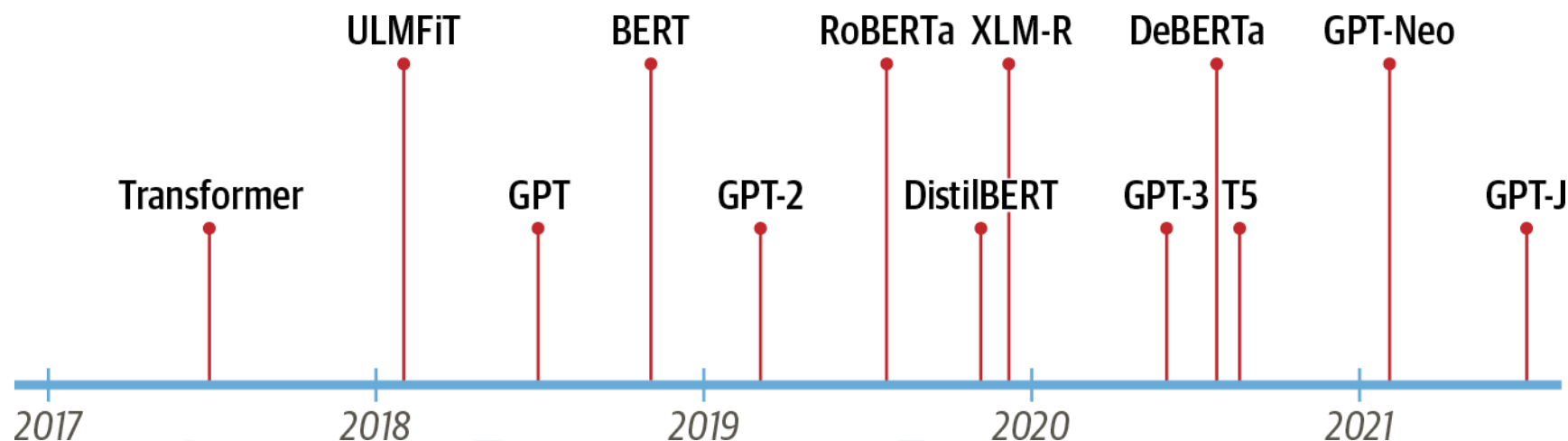


Figure 1-1. The transformers timeline

Agenda

- The Encoder-Decoder Framework
- Attention Mechanisms
- Transfer Learning in NLP
- Applications and Hugging Face Ecosystem

<https://ataghinezhad.github.io/>

The Encoder-Decoder Framework

- Purpose: Maps input sequences to output sequences
- Example: English “Transformers are great!” → German “Transformer sind grossartig!”
- Process:
 - Encoder: Converts input into a hidden state vector
 - Decoder: Generates output from this hidden state
- Reference: Sequence-to-sequence illustration (see Figure 1-3)

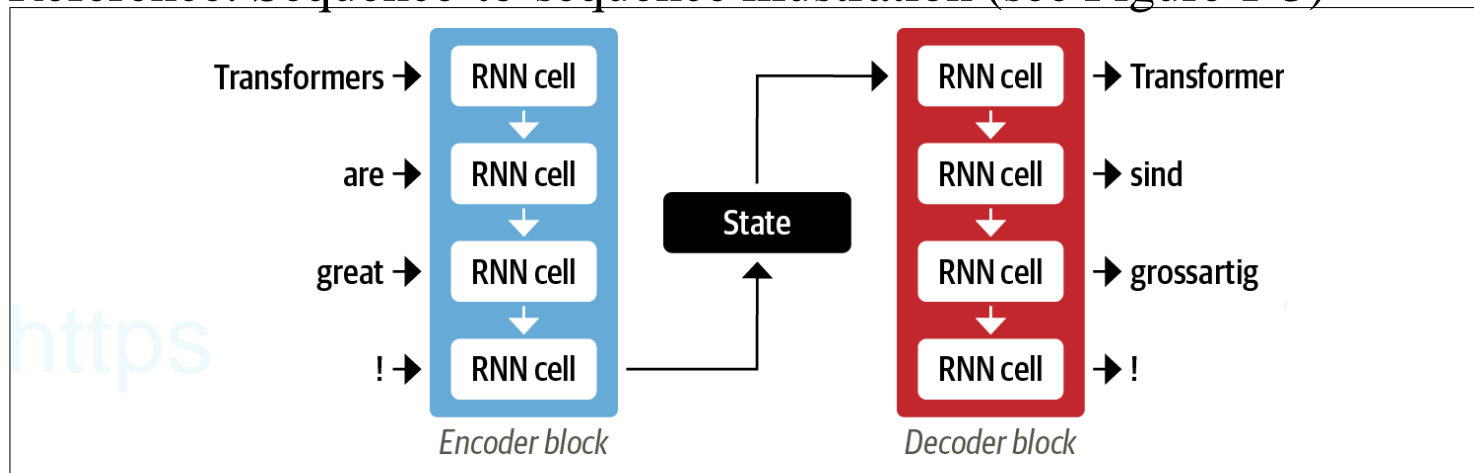


Figure 1-3. An encoder-decoder architecture with a pair of RNNs (in general, there are many more recurrent layers than those shown here)

Limitations of Encoder-Decoder Models

- Issue: Information Bottleneck
- Impact: Loss of detail, especially for long sequences
- Preview: Introduction of attention mechanisms to overcome these limitations
- Reference: Illustration of hidden state compression (see Figure 1-2)

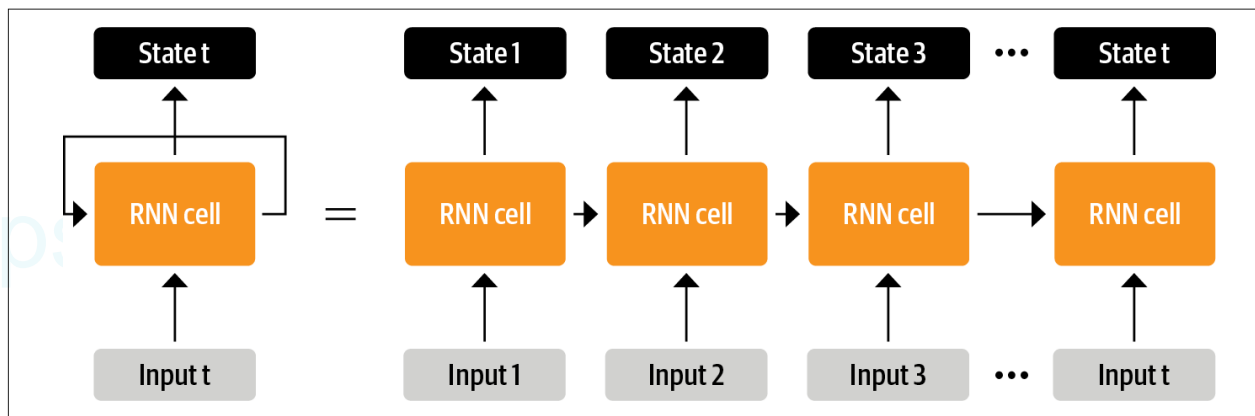


Figure 1-2. Unrolling an RNN in time

Attention Mechanisms

- Core Idea: Decoder accesses all encoder hidden states
- Mechanism: Assigns varying weights to different input tokens at each decoding step
- Benefit: Improves translation quality by focusing on the most relevant parts of the input
- Reference: Detailed attention process (see Figure 1-4)

<https://ataghinezhad.github.io/>

Visualizing Attention

- Example: English-to-French translation alignment
- Insight: Highlights the model's ability to align words between languages
- Reference: Attention weight visualization (see Figure 1-5)

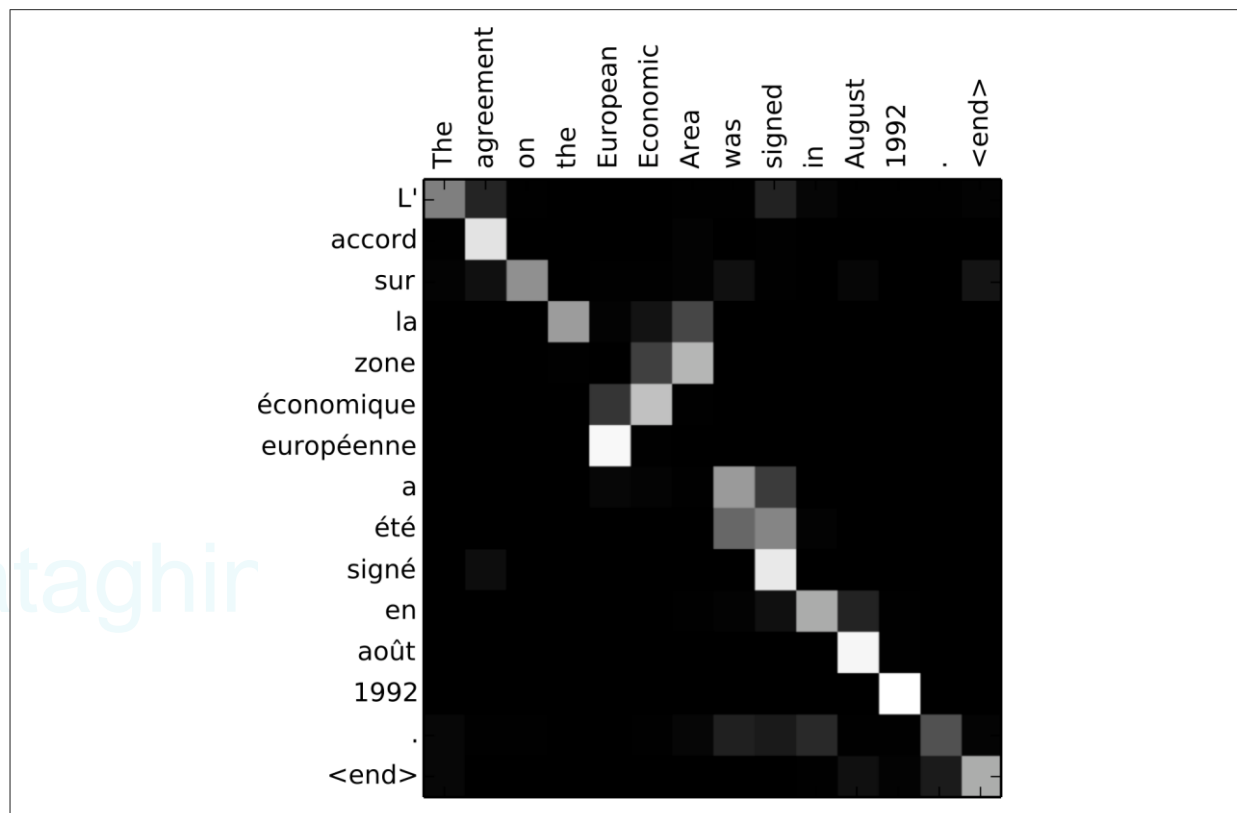


Figure 1-5. RNN encoder-decoder alignment of words in English and the generated translation in French (courtesy of Google)

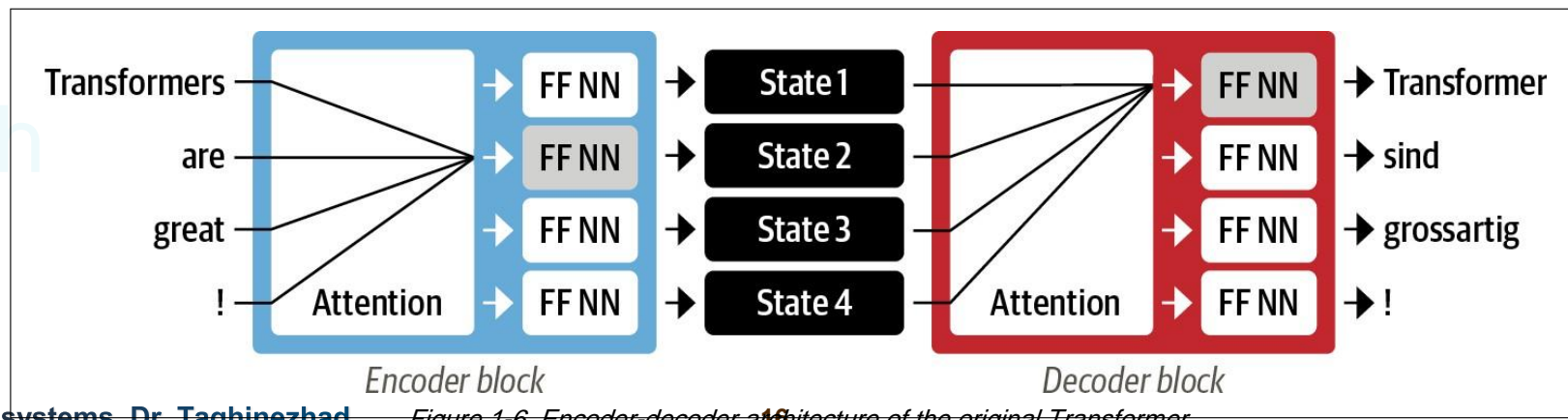
Limitations of RNN-Based Attention Models

- Challenge: Inherently sequential computations in RNNs
- Consequence: Reduced ability to parallelize processing
- Implication: Necessitated development of new architectures for efficiency

<https://ataghinezhad.github.io/>

Introduction to Transformers

- Innovation: Removal of recurrence; adoption of self-attention
- Key Feature: Self-attention operates on all sequence positions simultaneously
- Advantage: Enables parallel processing and faster training
- Reference: Transformer architecture overview (see Figure 1-6)
- Components:
 - Encoder with self-attention layers
 - Decoder with self-attention and feed-forward networks
- Outcome: Enhanced performance and scalability for NLP tasks
- Reference: Detailed transformer layout (see Figure 1-6)



Transfer Learning in NLP

- Background: Well-established in computer vision (e.g., ResNet + ImageNet)
- Challenge: Establishing a pretraining process in NLP
- Breakthrough: GPT and BERT leverage unsupervised learning for pretraining

<https://ataghinezhad.github.io/>

Transfer Learning Process

- Pretraining: Exposure to large, diverse corpora to learn general language features
- Fine-Tuning: Adapting the pretrained model to specific tasks with minimal labeled data
- Result: High-quality performance on downstream NLP tasks

<https://ataghinezhad.github.io/>

Impact of GPT and BERT

- **Eliminated Need:** No need for task-specific architectures built from scratch
- **Performance:** Consistently break NLP benchmarks
- **Emergence:** Catalyst for a diverse array of transformer-based models

<https://ataghinezhad.github.io/>

Hugging Face Ecosystem

- Tools and Libraries: Comprehensive suite for training and deployment of transformers
- Capabilities:
- Access to pre-trained models
- Simplified fine-tuning processes
- Community support and ongoing innovation

<https://ataghinezhad.github.io/>

The Emergence of Transfer Learning in NLP

- Timeframe: 2017–2018
- Key Insight: Unsupervised pretraining improves performance (OpenAI's sentiment analysis breakthrough)
- Milestone: Introduction of ULMFiT
- Reference: (See textual context, no specific figure)

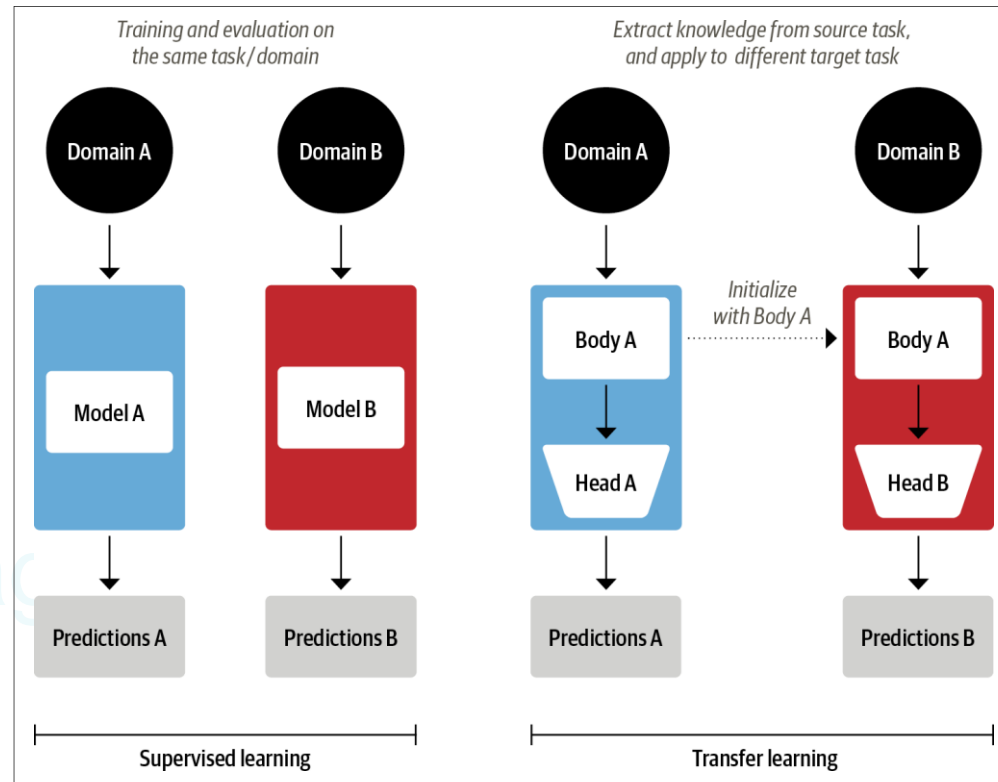


Figure 1-7. Comparison of traditional supervised learning (left) and transfer learning (right) In computer

The ULMFiT Process

- Three Main Steps (See Figure 1-8):
 - Pretraining: Learn language modeling on large, unlabeled corpora (e.g., Wikipedia)
 - Domain Adaptation: Adapt the model to in-domain data (e.g., from Wikipedia to IMDb reviews)
 - Fine-Tuning: Add and train a task-specific classification layer
- Reference: Figure 1-8 (courtesy of Jeremy Howard)

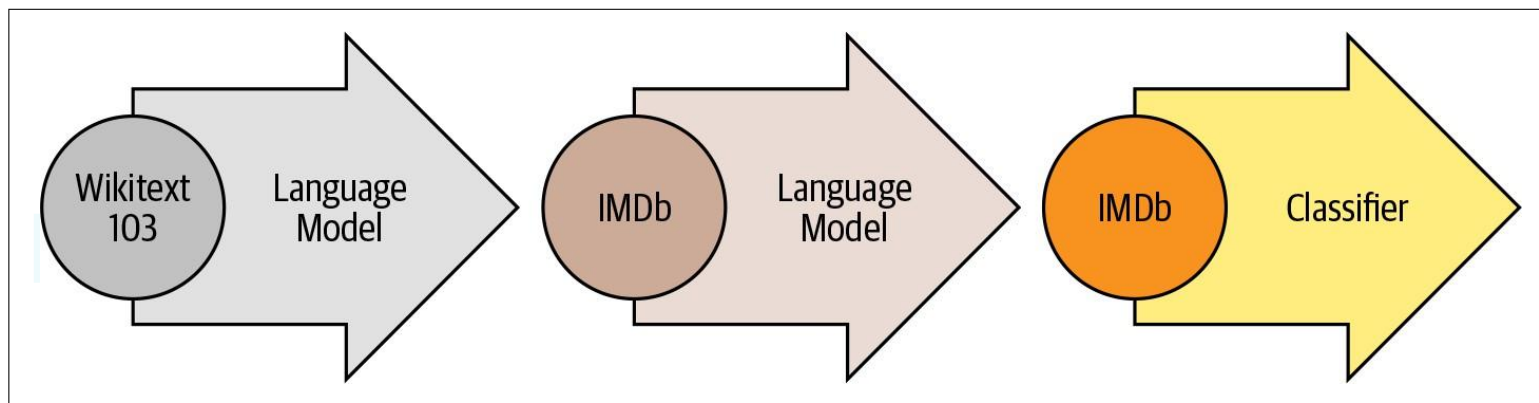


Figure 1-8. The ULMFiT process (courtesy of Jeremy Howard)

From ULMFiT to Transformers: GPT & BERT

- Building on ULMFiT: Transfer learning sets the stage for transformer models
- GPT:
 - Uses the decoder part of the Transformer
 - Pretrained on BookCorpus (7,000 unpublished books)
- BERT:
 - Uses the encoder part with masked language modeling
 - Pretrained on BookCorpus and English Wikipedia

<https://ataghinezhad.github.io/>

Bridging the Gap with Hugging Face Transformers

- Challenges:
 - Diverse frameworks (PyTorch, TensorFlow) and non-standardized implementations
- Solution: The Transformers library by Hugging Face
 - Provides a unified API for over 50 transformer architectures
 - Supports PyTorch, TensorFlow, and JAX
 - Offers task-specific heads for seamless fine-tuning
- Benefit: Accelerates integration from weeks to an afternoon

<https://ataghinezhad.github.io/>

A Tour of Transformer Applications

- Wide-Ranging Applications:
 - Text Classification
 - Named Entity Recognition (NER)
 - Question Answering
 - Summarization
 - Translation
 - Text Generation
- Note: Upcoming slides will include code examples and outputs for these tasks.

<https://ataghinezhad.github.io/>

Text Classification with Transformers

- Pipeline Example: `pipeline("text-classification")`
- Highlights:
 - Automatic download of pretrained weights
 - Supports sentiment analysis, multiclass, and multilabel classification
 - Outputs include confidence scores
- Reference: (See example code and outputs)

<https://ataghinezhad.github.io/>

Named Entity Recognition (NER)

- Pipeline Example: `pipeline("ner", aggregation_strategy="simple")`
- Functionality:
 - Detects entities like organizations, locations, and persons
 - Groups multi-word entities (e.g., “Optimus Prime”) into single units
- Reference: (Refer to example code and output)

<https://ataghinezhad.github.io/>

Question Answering

- Pipeline Example: `pipeline("question-answering")`
- Process:
 - Input: A passage of text (context) and a specific question
 - Output: The answer span with character indices (start and end)
- Example: Extracting “an exchange of Megatron” from customer feedback
- Reference: (Refer to example code and output)

<https://ataghinezhad.github.io/>

Summarization

- Pipeline Example: `pipeline("summarization")`
- Key Parameters:
 - `max_length` controls summary length
 - `clean_up_tokenization_spaces` ensures proper formatting
- Example: Generates a concise summary of a customer complaint
- Reference: (Refer to example code and output)

<https://ataghinezhad.github.io/>

Translation

- Pipeline Example:
`pipeline("translation_en_to_de", model="Helsinki-NLP/opus-mt-en-de")`
- Features:
 - Overrides the default model for targeted performance
 - Produces translations with correct formal language usage
- Example: Translating customer feedback from English to German
- Reference: (Refer to example code and output)

<https://ataghinezhad.github.io/>

Text Generation

- Pipeline Example: `pipeline("text-generation")`
- Use Case: Autocompleting or generating extended customer service responses
- Process:
 - Combine an input text with a prompt to generate a response
- Example: Generating a customer service reply
- Reference: (Refer to example code and output)

<https://ataghinezhad.github.io/>

The Hugging Face Ecosystem

- Key Components:
 - A family of libraries (e.g., Transformers)
 - The Hugging Face Hub for pretrained models, datasets, and evaluation scripts
- Visual Reference: See Figure 1-9 (Ecosystem Diagram)
- Benefit: Accelerates research, reproducibility, and deployment

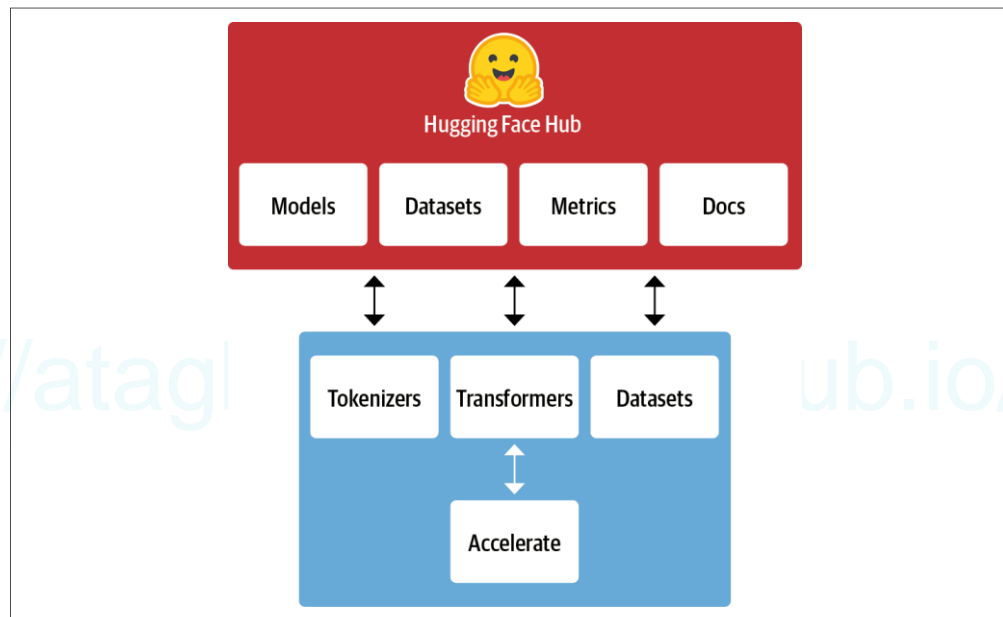


Figure 1-9. An overview of the Hugging Face ecosystem

The Hugging Face Hub in Detail

- Hub Features:
 - Over 20,000 freely available models
 - Filters for task, framework, and dataset (see Figures 1-10)
 - Interactive widgets for trying models directly (see Figure 1-11)
 - Detailed model and dataset cards for documentation
- Reference: Figures 1-10 and 1-11

The image shows two side-by-side screenshots from the Hugging Face Hub. The left screenshot displays the 'Models' page with a search bar, navigation tabs (Models, Datasets, Resources, Solutions, Pricing), and a list of models. The right screenshot shows a detailed 'Model card' for the 'BERT base model (uncased)', including a description, disclaimer, and model description.

Models Page (Left Screenshot):

- Search: Search models, datasets, users...
- Navigation: Models, Datasets, Resources, Solutions, Pricing
- Tasks: Fill-Mask, Question Answering, Summarization, Table Question Answering, Text Classification, Text Generation, Text2Text Generation, Token Classification, Translation, Zero-Shot Classification, Sentence Similarity (+12)
- Libraries: PyTorch, TensorFlow, JAX (+19)
- Datasets: wikipedia, common_voice, bookcorpus, dcep europarl-jrc-acquis, glue, squad
- Models List (25,493):
 - bert-base-uncased (Fill-Mask, Updated May 18, 27.5M, 42)
 - xlm-roberta-base (Fill-Mask, Updated Sep 16, 5.88M, 9)
 - roberta-large (Fill-Mask, Updated May 21, 5.26M, 15)
 - distilbert-base-uncased (Fill-Mask, Updated Aug 29, 4.86M, 22)
 - gpt2 (Text Generation, Updated May 19, 4.64M, 15)

Model Card (Right Screenshot):

- Model card: Files, Settings, Train, Deploy, Use in Transformers
- Model: BERT base model (uncased)
- Description: Pretrained model on English language using a masked language modeling (MLM) objective. It was introduced in [this paper](#) and first released in [this repository](#). This model is uncased: it does not make a difference between english and English.
- Disclaimer: The team releasing BERT did not write a model card for this model so this model card has been written by the Hugging Face team.
- Model description: BERT is a transformers model pretrained on a large corpus of English data in a self-supervised fashion. This means it was pretrained on the raw texts only, with no humans labelling them in any way (which is...)
- Hosted inference API: Select AutoNLP in the "Train" menu to fine-tune this model automatically. Downloads last month: 27,529,242.
- Fill-Mask: Mask token: [MASK]. Input: I looked at my [MASK] and saw I was la. Compute button.
- Computation time on cpu: cached. Performance metrics: watch (0.228), phone (0.099), face (0.059), hands (0.053), hand (0.035).
- JSON Output, Maximize

Figure 1-10. The Models page of the Hugging Face Hub, showing filters on the left and a list of models on the right

Hugging Face Tokenizers

- Tokenization in NLP:
 - Splits raw text into tokens (words, subwords, characters)
 - Essential for converting text into numerical representations for transformers
- Hugging Face Tokenizers:
 - Provides various tokenization strategies
 - Extremely fast due to its Rust backend
 - Manages pre- and post-processing (normalization, formatting)
 - Loads tokenizers just like pretrained model weights in Transformers

<https://ataghinezhad.github.io/>

Hugging Face Datasets

- Challenges in Data Handling:
 - Loading, processing, and storing large datasets
 - Custom scripts often needed for data download and format conversion
- Hugging Face Datasets:
 - Standard interface for thousands of datasets from the Hub
 - Smart caching to avoid repeated preprocessing
 - Memory mapping to overcome RAM limitations
 - Seamless integration with Pandas and NumPy
 - Provides metric scripts for reproducible and trustworthy evaluations

<https://atahinezhad.github.io/>

Hugging Face Accelerate

- Challenges with Custom Training Scripts:
 - Porting code from a laptop to a large-scale cluster can be problematic
- Hugging Face Accelerate:
 - Adds abstraction to training loops in PyTorch and other frameworks
 - Simplifies infrastructure changes for distributed training
 - Accelerates workflow by managing custom logic behind the scenes

<https://ataghinezhad.github.io/>

Main Challenges with Transformers

- Language:
 - Research is predominantly focused on English
 - Scarcity of pretrained models for rare or low-resource languages
- Data Availability:
 - Despite transfer learning, large amounts of labeled data are still needed
- Working with Long Documents:
 - Self-attention is effective for paragraphs but becomes expensive for longer texts
- Opacity:
 - Transformer models are complex and often lack interpretability
- Bias:
 - Pretraining on internet text can imprint undesirable biases (racist, sexist, etc.)

<https://ataghinezhad.github.io/>

Conclusion & Looking Ahead

- Recap:
 - Tokenizers, Datasets, and Accelerate form the backbone of Hugging Face's ecosystem
 - Each component plays a crucial role in training, evaluating, and deploying transformer models
- Challenges:
 - Addressing language limitations, data needs, document length, opacity, and bias
- Next Steps:
 - Upcoming chapters will provide hands-on experience with text classification and other applications

<https://ataqinezhad.github.io/>

End of Chapter 1

<https://ataghinezhad.github.io/>