



# Machine Learning Lifecycle and Standards for Transparency and Accountability

Christopher Klamm

@chklamm

24.10.2022

**“It seemed so simple: share all data, code and parameter settings, and other researchers will be able to obtain the same results.”** [Belz et al. 2021](#)



**Nils Reimers** @Nils\_Reimers · 4 Std.

Antwort an @PreetumNakkiran

Same. Arxiv + Tweet + reproducible results + actually usable code for other use cases makes much more sense than the hassle and limitations of conferences.

*Director of Machine Learning at cohore.ai*

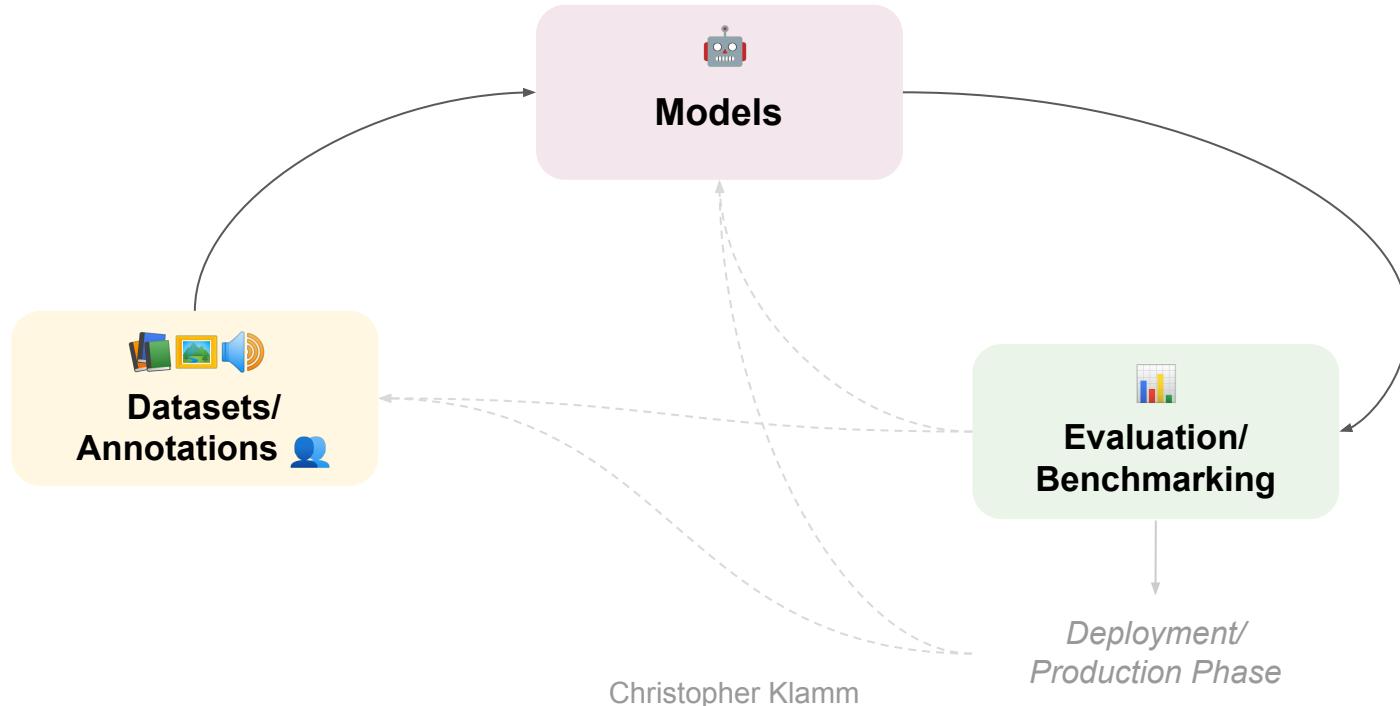
23.10.2022

...

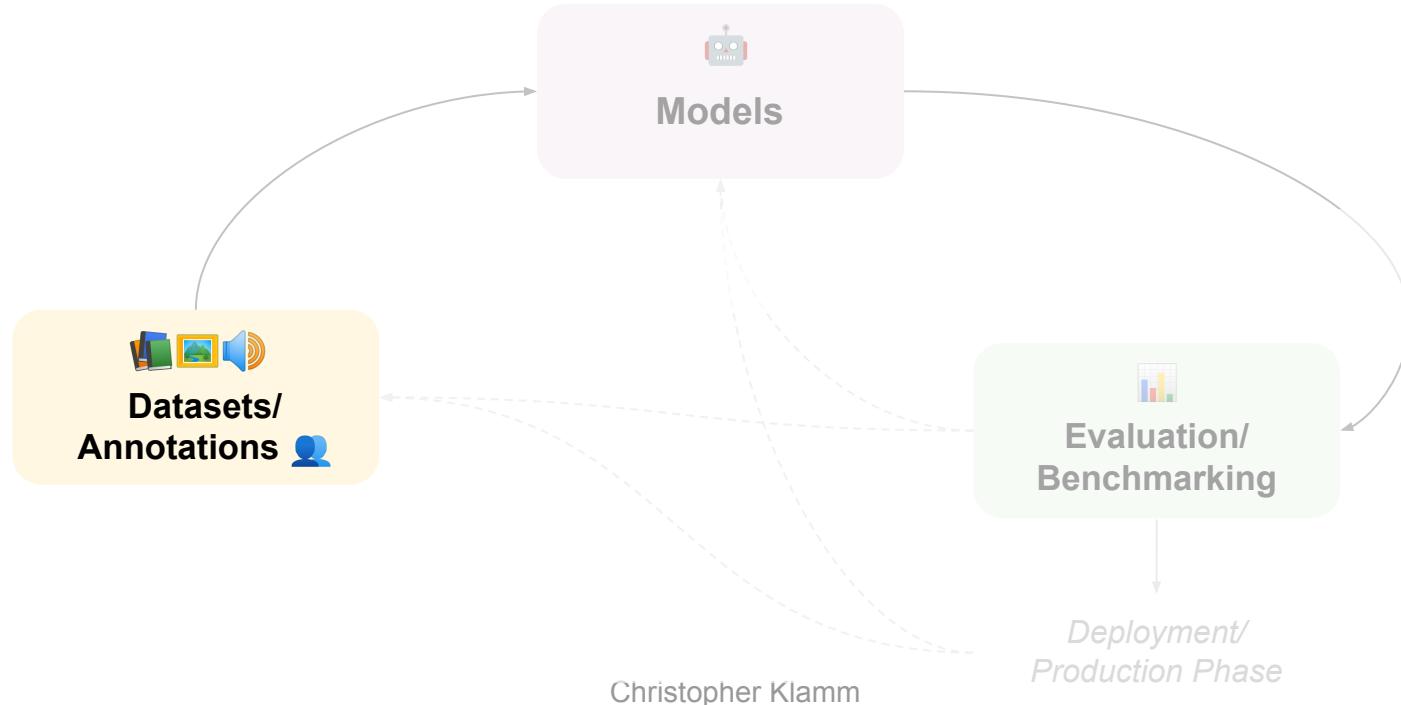
**BUT**

“ [...] a tiny **14.03%** of the **513** original/reproduction score pairs we looked at were the same. [...] [S]mall differences in code have been found to result in big differences in performance.” [Belz et al. 2021](#)

# Machine Learning Lifecycle/ Pipeline



# Machine Learning Lifecycle/ Pipeline



# Datasets

# Dataset

*A dataset is a collection of [data type] combined with annotations.*

## Text

e.g., news, blogs, tweets, comments, etc.

## Audio

e.g., audiobooks, speeches, podcasts, etc.

## Images/ Video

e.g., youtube/ tiktok/ instagram (clips), movies etc.

Source: <https://youtu.be/lim-a-dmND8>



*speech w/ captions (audio, video & text)*

... or **multimodal** datasets combining text, images, videos or other additional signals.

# Machine Learning/ NLP Tasks

- **Machine Translation** (e.g., EN-DE, EN-SE, ...)
- **Summarisation** (e.g., tl;dr)
- **Question & Answering** (e.g., chatbots)
- **Speech-to-Text** (e.g., caption generation)
- **Classification** (e.g., hate speech detection)
- **Generation** (e.g., arguments)
- ...

# Machine Learning/ NLP Tasks

- Machine Translation
  - Summarisation
  - Question & Answer
  - Speech-to-Text
  - Classification
  - Generation

**Input:** I was late for my meeting because I started playing video games for two minutes later than I intended.

**Model Output:** The person was late for their meeting because they started playing video games for only playing for two minutes.

**PaLM** *Chowdhery et al. 2022: 38*

# Explaining a Joke

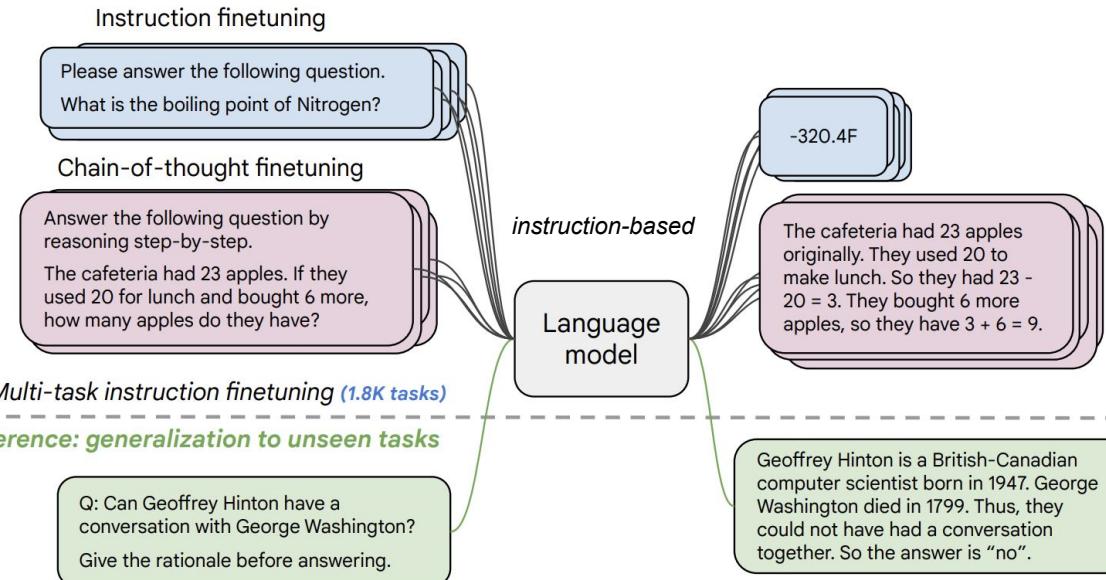
**Input:** I was supposed to start writing the paper at 5:00 PM. But then I started playing with this cool new language model for 10 minutes. 10 minutes later, it's suddenly 9:30 PM!

**Model Output:** This joke is about how time flies when you're having fun. The person was supposed to start writing a paper at 5:00 PM, but they started playing with a language model instead. They thought they were only playing with it for 10 minutes, but it was actually 4.5 hours.

# Machine Learning/ NLP Tasks

(20.10.22) **FLAN-T5** [Chung et al. 2022](#), [Demo](#)

- Machine
- Summar
- Question
- Speech-
- Classific
- Generat
- ...



# Machi

- Mac
- Sun
- Que
- Spe
- Clas
- Gen
- ...

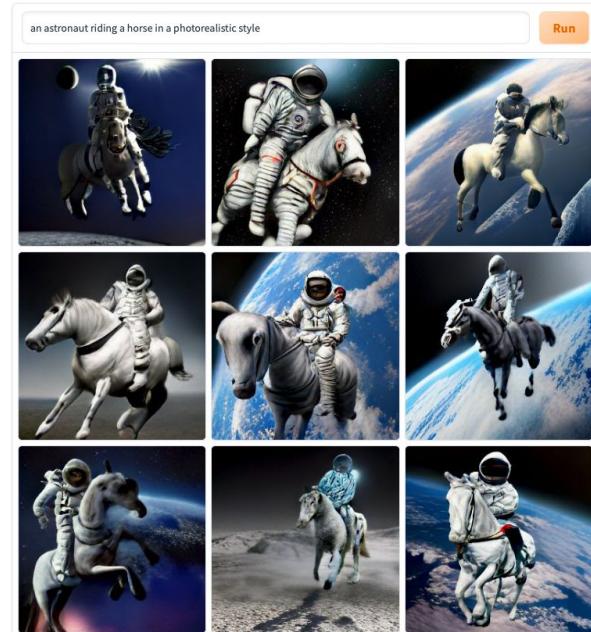
## TEXT DESCRIPTION

An astronaut Teddy bears A bowl of soup  
riding a horse lounging in a tropical resort in space playing basketball with cats in space  
in a photorealistic style in the style of Andy Warhol as a pencil drawing



## DALL·E mini by [craiyon.com](#)

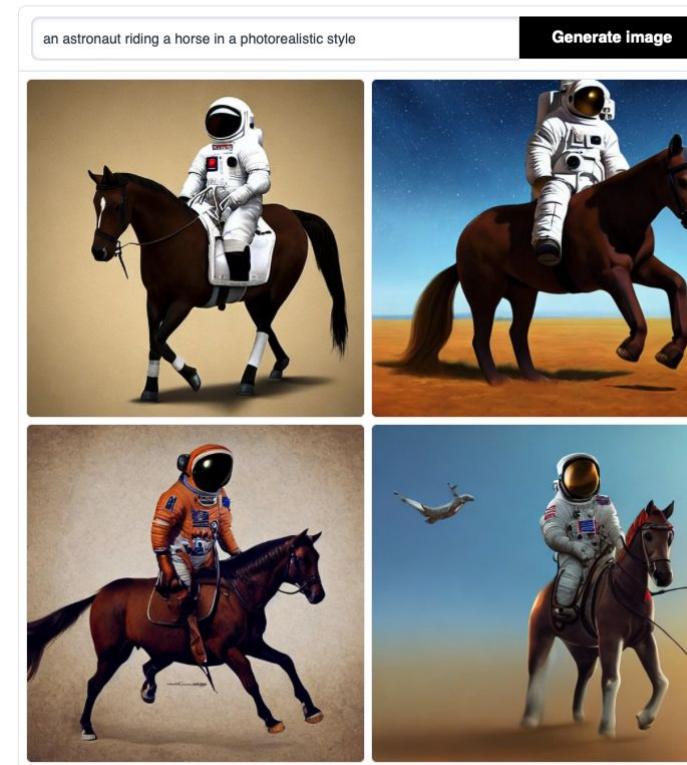
AI model generating images from any prompt!



<https://openai.com/dall-e-2/>, Demo

# Machine Learning/ NLP

- Machine Translation
- Summarisation
- Question & Answering
- Speech-to-Text
- Classification
- Generation
- ...



[Rombach et al. 2022, Demo](#)

# Annotations

*A dataset is a collection of [data type] combined with annotations.*

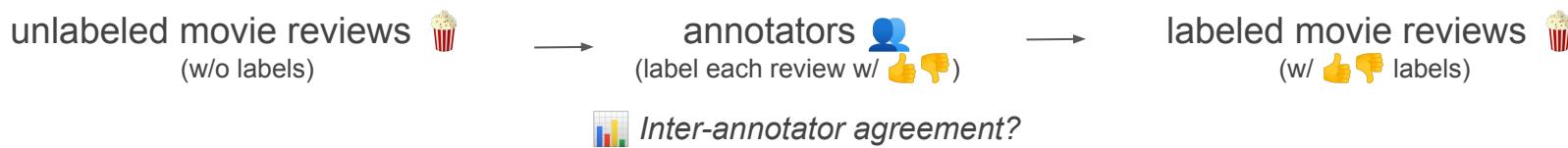
## Task *Sentiment Classification*

**input text**      This film is terrible. You don't really need to read this review further.

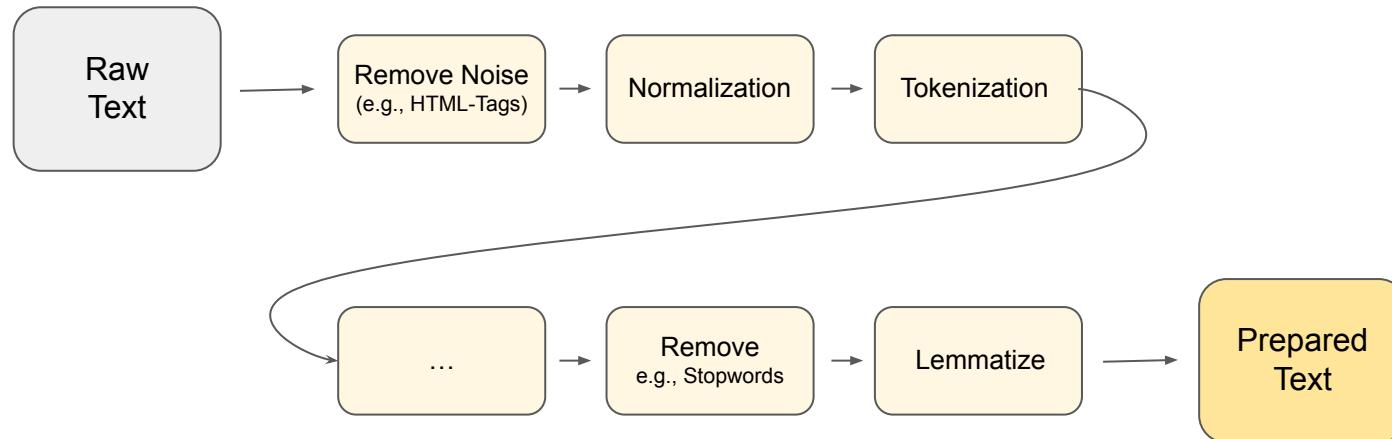
**labels**      positive      OR      negative

**annotation**      negative

## Annotation process



# Pre-Processing/ Cleaning



# Datasets

*Data Statements* ([Bender/ Friedman 2018](#)) & *Datasheets* ([Gebru et al. 2021](#))

**Data Statements:** “*A data statement is a characterization of a dataset that provides context to allow developers and users to better understand how experimental results might generalize, how software might be appropriately deployed, and what biases might be reflected in systems built on the software.*” ([Bender/ Friedman 2018: 1](#))

**Datasheets:** “*A datasheet is describing the operating characteristics, test results, recommended usage, and other information of a dataset. It documents the motivation, composition, collection process, recommended uses and so on.*” ([Gebru et al. 2021: 1-2](#))

# Datasheet in the PaLM (2022) paper I

**purpose**

<b>Motivation</b>	
For what purpose was the dataset created? Who created the dataset? Who funded the creation of the dataset?	The dataset was created for pre-training language models by a team of researchers at Google.
Any other comments?	<p>To train the model, we started with the dataset described in <a href="#">Du et al. (2021)</a>. The dataset is representative of a wide range of natural language use cases, and contains a high-quality filtered subset of webpages combined with books, Wikipedia pages, and data from public domain social media conversations used by <a href="#">Adiwardana et al. (2020)</a>. We made several modifications to the dataset:</p> <ul style="list-style-type: none"><li>• Adjusted the mixing proportions of the components of the dataset to avoid repeating training examples and minimize the risk of unstable training or overfitting. The mixing proportions are given in <a href="#">Table 2</a>.</li><li>• Used multilingual versions of Wikipedia and conversations data to improve the multilingual capabilities of the model and increase the number of tokens. The training mixture includes 124 languages, with English accounting for approximately 78% of the training tokens. The language distribution is shown in <a href="#">Figure 28</a>.</li><li>• Included deduplicated code from GitHub, filtered by license so as to exclude repositories with a copyleft license.</li><li>• Included date markers in the conversation data, to allow for conditional generation on dates (in particular conditioning on a recent date to avoid generating outdated facts).</li></ul>

# Datasheet in the PaLM (2022) paper II

## dataset type

Composition	
What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?	All instances of the dataset are text-only documents. Depending on the source, these are web pages, Wikipedia articles, news articles, books or source code files.
How many instances are there in total (of each type, if appropriate)?	The data makeup is given in Table 2.
Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?	The dataset is a (random) sample from a larger set. The sampling methodology is described in <a href="#">Du et al. (2021)</a> . The different components of the dataset have different weights, as specified in Table 2.
What data does each instance consist of?	Each instance is a SentencePiece ( <a href="#">Kudo &amp; Richardson, 2018b</a> ) encoded sequence of text.
Is there a label or target associated with each instance?	No, there are no labels associated with each instance.
Is any information missing from individual instances?	No.
Are relationships between individual instances made explicit?	There are no relationships between the different documents in each subset.
Are there recommended data splits?	We use random splits for the training and development sets.
Are there any errors, sources of noise, or redundancies in the dataset?	Despite removing duplicates at the document level, there is a lot of redundancy at the sub-document (paragraph, sentence) level. There is also redundancy coming from different instantiations of the same textual pattern, and from general low quality text from the Web, e.g., SEO spam.
Is the dataset self-contained, or does it link to or otherwise rely on external resources?	The dataset is self-contained.
Does the dataset contain data that might be considered confidential?	No.
Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?	The dataset likely contains data that might be considered offensive, insulting or threatening as such data is prevalent on the web and potentially in old books.

# Datasheet in the PaLM (2022) paper III

Collection Process	
How was the data associated with each instance acquired?	The data was collected from publicly available sources.
What mechanisms or procedures were used to collect the data?	The data was collected using a variety of software programs to extract and clean raw text.
If the dataset is a sample from a larger set, what was the sampling strategy?	<p>The sampling methodology is described in <a href="#">Du et al. (2021)</a>. For Web documents, we use two methods of sampling:</p> <ul style="list-style-type: none"> <li>• Random sampling based on a classifier that gives higher probability to high quality documents.</li> <li>• Selecting documents that are also in the Colossal Clean Crawled Corpus (C4) (<a href="#">Raffel et al., 2020</a>).</li> </ul>
Who was involved in the data collection process?	A team of researchers at Google.
Over what timeframe was the data collected?	2019-2021
Were any ethical review processes conducted?	No.
Preprocessing, cleaning, and labeling	
Was any preprocessing, cleaning, or labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)?	We remove boilerplate from web pages using proprietary software. We also remove HTML markup. We extract conversations using a special-purpose algorithm.
Is the software used to preprocess, clean, or label the instances available?	No.

pre-processing

# Datasheet in the PaLM (2022) paper IV

Uses	
Has the dataset been used for any tasks already?	Yes, we use the dataset for pre-training language models.
Is there a repository that links to any or all papers or systems that use the dataset?	No.
What (other) tasks could the dataset be used for?	The large-scale task-agnostic nature of the dataset makes it suitable for many NLP tasks such as language model pretraining or question answering.
Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?	The dataset is static in nature and thus will become progressively more "stale". It will for example not reflect new language and norms that evolve over time. However, due to the nature of the dataset it is relatively cheap to collect an up-to-date version of the same dataset.
Are there tasks for which the dataset should not be used?	This model should not be used for any of the unacceptable language model use cases, e.g., generation of toxic speech.
Distribution	
Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?	No.

unintended use

# Creating Datasheets is a challenging task

Heger et al. 2022 evaluated Datasheets in a study with 14 ML practitioners

- “One of our most concerning findings is that participants did not make the connection between the datasheets for datasets questions and their responsible AI implications [...]” (p. 22)  
→ we need to find out “How best to train ML researchers and practitioners to [...] **anticipate potential consequences of their work** is an area where more research is needed.” (p. 23)
- Other **improvements for Datasheets** (p. 23ff):
  - “Make explicit the connection between data documentation and responsible AI”
  - “Make data documentation frameworks practical”
  - “Adapt data documentation frameworks to different contexts”
  - “Don’t automate away responsibility, but do support simple tasks with automation”
  - “Clarify the target audience for data documentation”
  - “Standardize and centralize data documentation”
  - “Integrate data documentation frameworks into existing tools and workflows”

see Heger et al. 2022

# Responsible Data Use Checklist (Rogers/ Baldwin/ Liens 2021)

## For papers using a previously-published resource:

1.  The authors explain their choice of data, given the available resources and their known limitations (e.g. representativeness issues, biases, annotation artifacts) and any data protection issues (e.g. inclusion of sensitive health data). *See Section* ---
2.  The authors discuss whether their use of a previously-published resource is compatible with its original purpose and license, and any known limitations (e.g. if the target user group is represented in the sample). *See Section* ---

## For papers contributing a new resource:

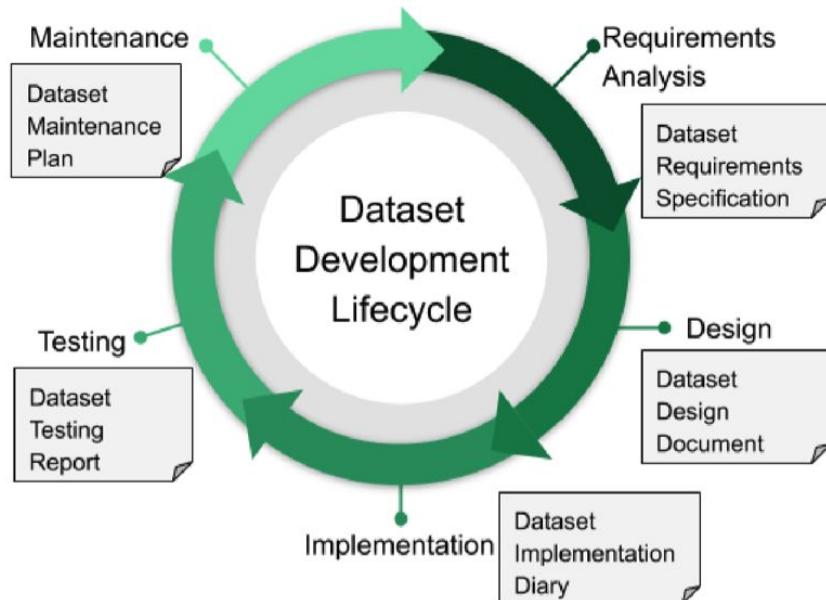
1.  The authors have the **legal basis** for processing the data and, if it is made public, for distributing it. (*Check one*)
  - 1.1.  The data is in public domain, and licensed for research purposes;
  - 1.2.  The data is used with consent of its creators or copyright holders;
  - 1.3.  If the data is used without consent, the paper makes the case to justify its legal basis (e.g. research performed in the public interest under GDPR). *See Section* ---
2.  The paper describes in detail the **full data collection protocol**, including collection, annotation, pre-processing, and filtering procedures. In the case that the dataset involves work by human subjects (e.g. data creation or annotation), the paper describes efforts to ensure fair compensation. *See Section* ---

## 3. Safe use of data is ensured. (*Check all that apply*)

- 3.1.  The data does not include any protected information (e.g. sexual orientation or political views under GDPR), or a specified exception applies. *See Section* ---
  - 3.2.  The paper is accompanied by a data statement describing the basic demographic and geographic characteristics of the population that is the source of the language data, and the population that it is intended to represent. *See Section* ---
  - 3.3.  If applicable: the paper describes whether any characteristics of the human subjects were self-reported (preferably) or inferred (in what way), justifying the methodology and choice of description categories. *See Section* ---
  - 3.4.  The paper discusses the harms that may ensue from the limitations of the data collection methodology, especially concerning marginalized/vulnerable populations, and specifies the scope within which the data can be used safely. *See Section* ---
  - 3.5.  If any personal data is used: the paper specifies the standards applied for its storage and processing, and any anonymization efforts. *See Section* ---
  - 3.6.  If the individual speakers remain identifiable via search: the paper discusses possible harms from misuse of this data, and their mitigation. *See Section* ---
4.  If any data or models are made public: **safe reuse** is ensured. (*Check all that apply*)
    - 4.1.  The data and/or pretrained models are released under a specified license that is compatible with the conditions under which access to data was granted (in particular, derivatives of data accessed for research purposes should *not* be deployed in the real world as anything other than a research prototype, especially commercially). *See Section* ---
    - 4.2.  The paper specifies the efforts to limit the potential use to circumstances in which the data/models could be used safely (such as an accompanying data/model statements). *See Section* ---
  5.  The data collection protocol was **approved by the ethics review board** at the authors' institution, or such review is not applicable for specified reasons. *See Section* ---

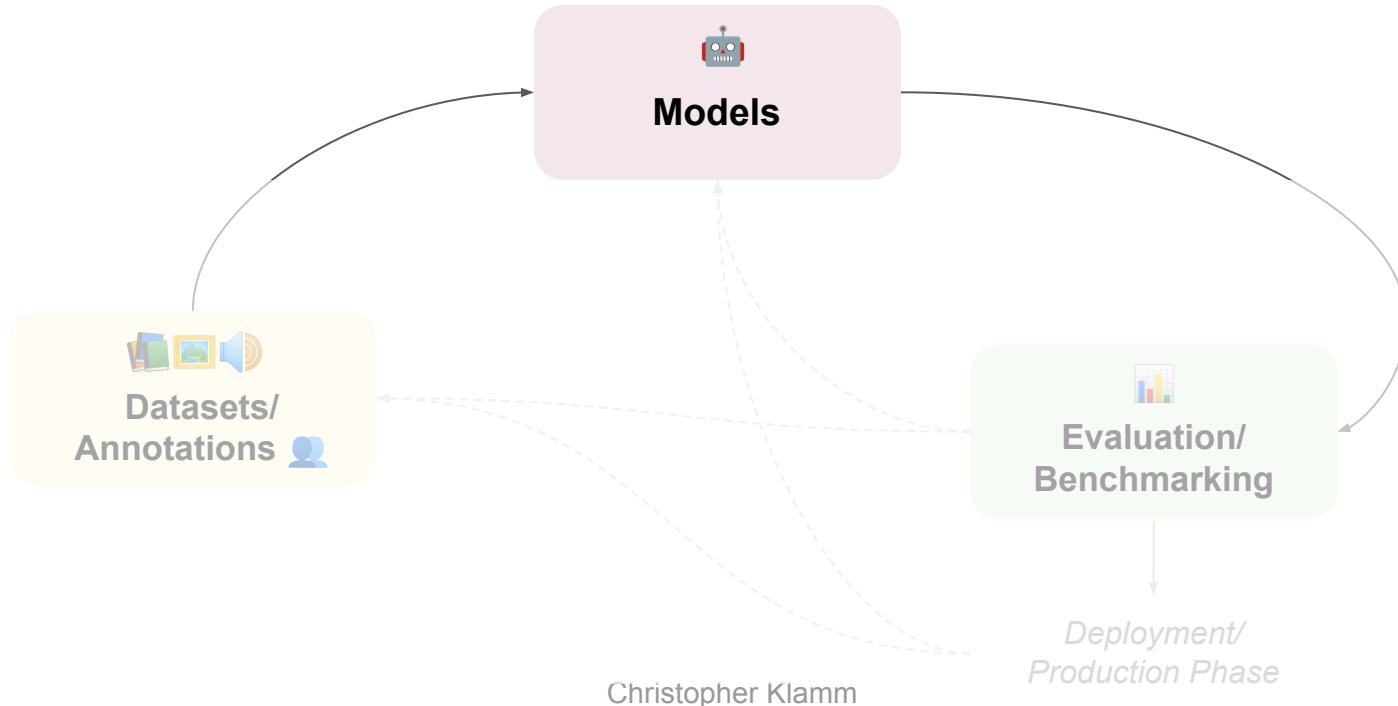
Rogers/ Baldwin/ Liens 2021: 4827

# Dataset Development Lifecycle (Hutchinson et al. 2020)



<i>Requirements analysis</i>	Deliberations about intentions, consultations with stakeholders, and analysis of use cases determine what data is required.
<i>Design</i>	Research is performed and subject matter experts are consulted in order to determine whether the data requirements can be met, and if so how best to do so.
<i>Implementation</i>	Design decisions are transformed into technologies such as software systems, annotator guidelines, and labeling platforms. Actions may employing and managing teams of human expert raters.
<i>Testing</i>	Data is evaluated and decisions about whether or not to use it are made.
<i>Maintenance</i>	Once collected, a dataset requires a large set of affordances, including tools, policies and designated owners.

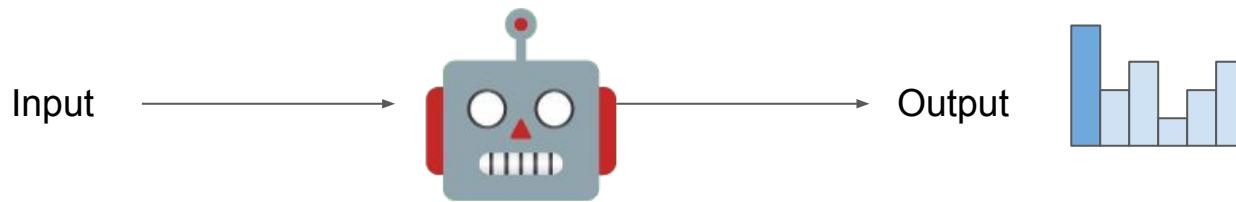
# Machine Learning Lifecycle/ Pipeline



# Models

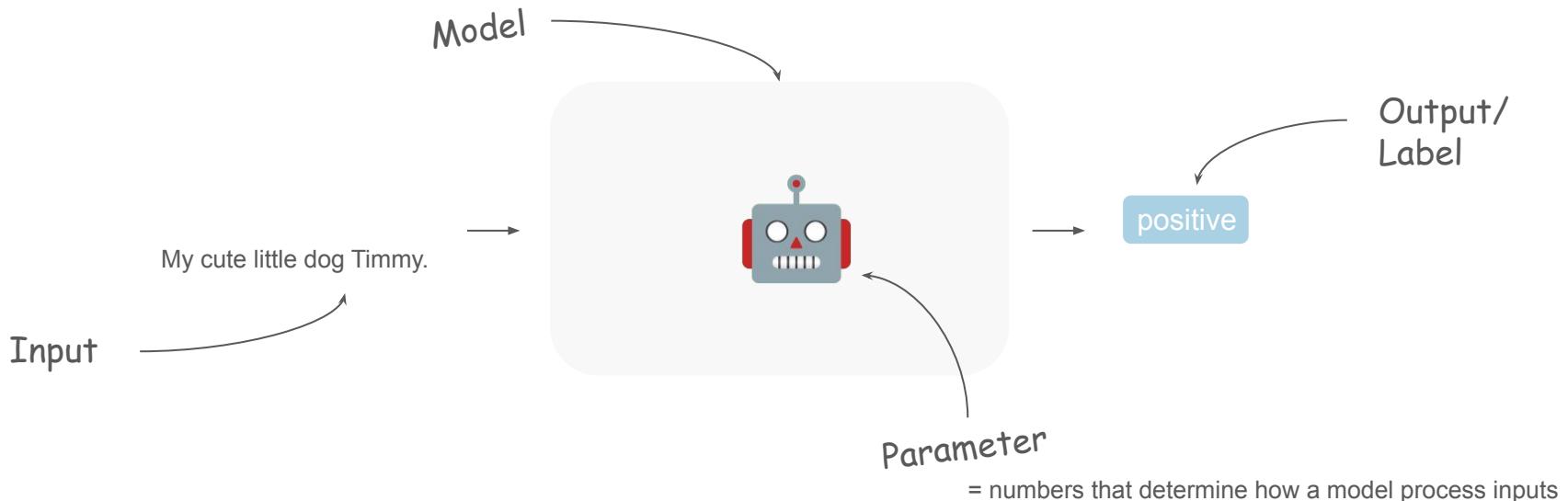


# Machine Learning



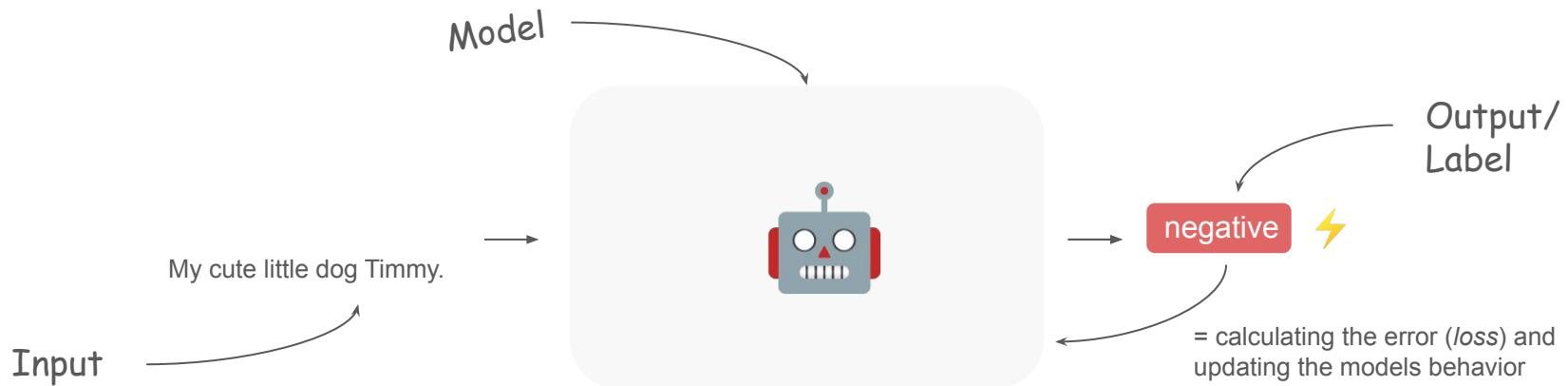


# Machine Learning



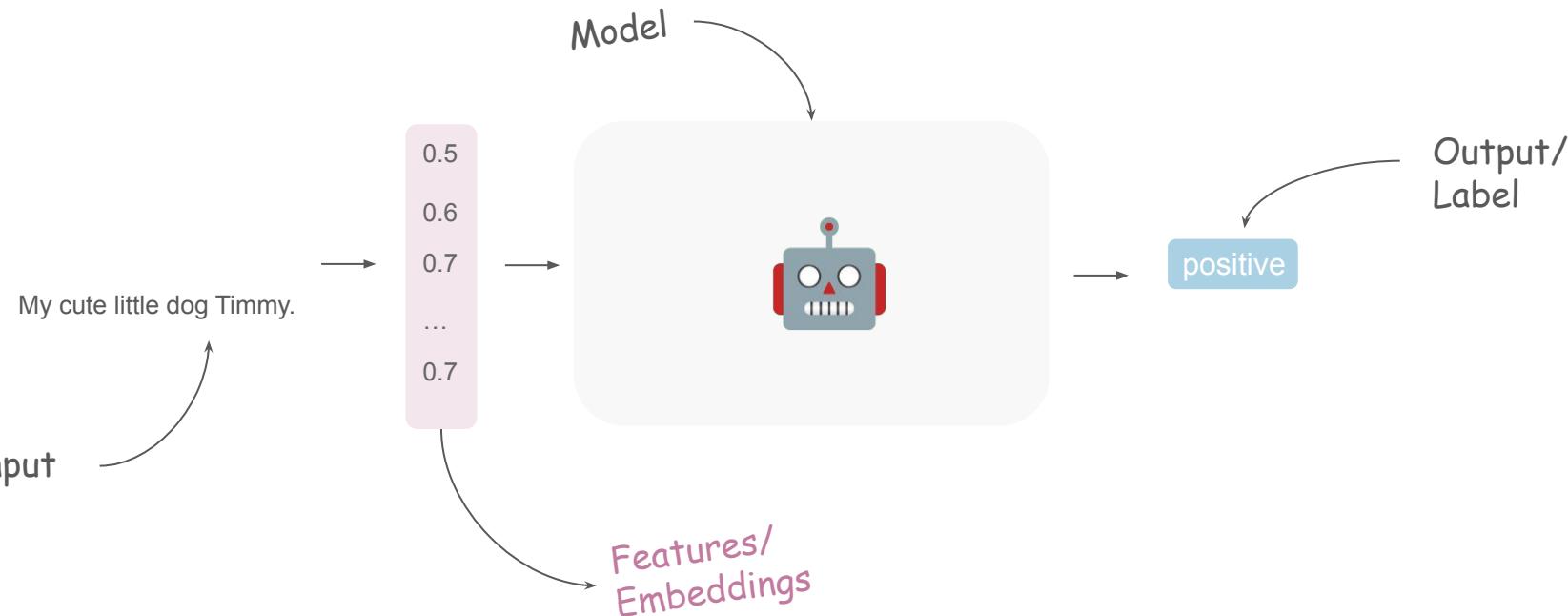


# Machine Learning



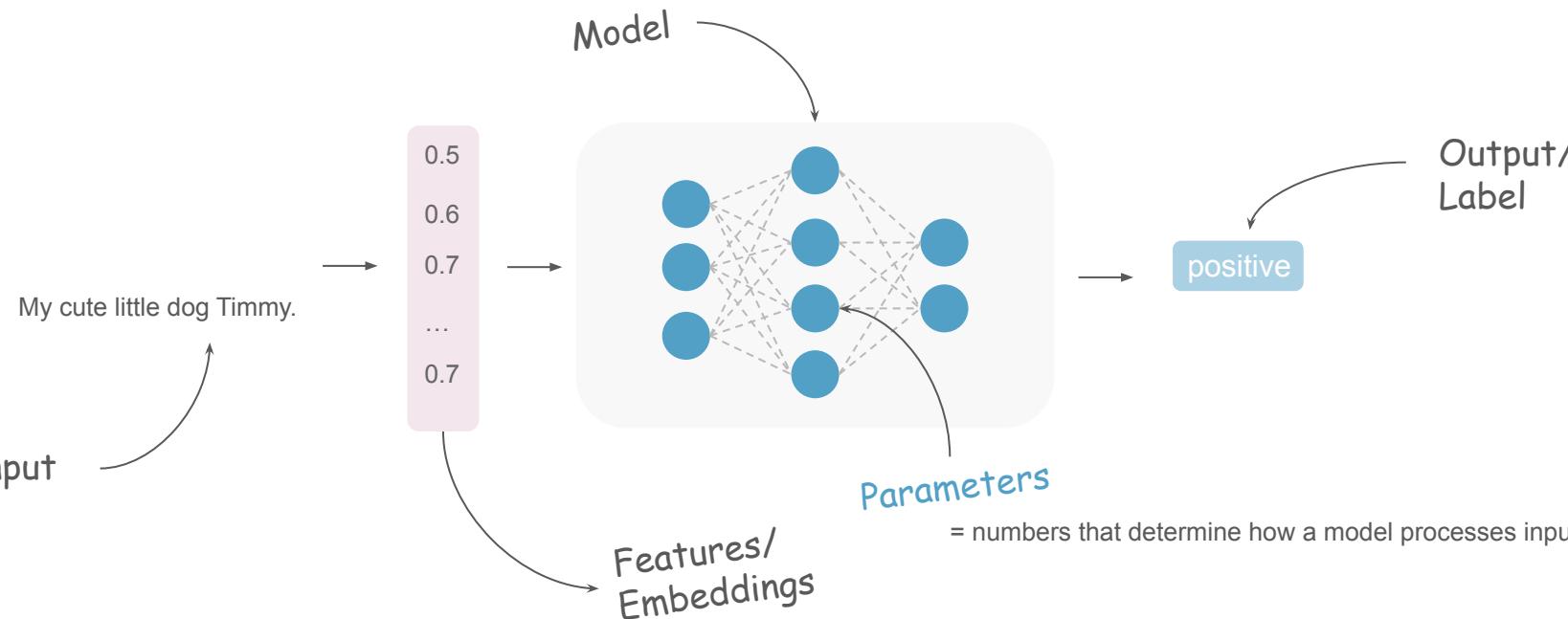


# Machine Learning





# Machine Learning



# Model Cards ([Mitchell et al. 2019](#))

## Model Card

- **Model Details.** Basic information about the model.
  - Person or organization developing model
  - Model date
  - Model version
  - Model type
  - Information about training algorithms, parameters, fairness constraints or other applied approaches, and features
  - Paper or other resource for more information
  - Citation details
  - License
  - Where to send questions or comments about the model
- **Intended Use.** Use cases that were envisioned during development.
  - Primary intended uses
  - Primary intended users
  - Out-of-scope use cases
- **Factors.** Factors could include demographic or phenotypic groups, environmental conditions, technical attributes, or others listed in Section 4.3.
  - Relevant factors
  - Evaluation factors
- **Metrics.** Metrics should be chosen to reflect potential real-world impacts of the model.
  - Model performance measures
  - Decision thresholds
  - Variation approaches
- **Evaluation Data.** Details on the dataset(s) used for the quantitative analyses in the card.
  - Datasets
  - Motivation
  - Preprocessing
- **Training Data.** May not be possible to provide in practice. When possible, this section should mirror Evaluation Data. If such detail is not possible, minimal allowable information should be provided here, such as details of the distribution over various factors in the training datasets.
- **Quantitative Analyses**
  - Unitary results
  - Intersectional results
- **Ethical Considerations**
- **Caveats and Recommendations**

# Model Card in the PaLM (2022) paper I

Example



Model Summary	
Model Architecture	Dense decoder-only model with 540 billion parameters. Transformer model architecture with variants to speed up training and inference. For details, see Model Architecture (Section 2).
Input(s)	The model takes text as input.
Output(s)	The model generates text as output.
Usage	
Application	<p>The primary use is research on language models, including: research on NLP applications like machine translation and question answering, advancing fairness and safety research, and understanding limitations of current LLMs.</p> <p>Within Google, PaLM is being used for research on a variety of open-ended text and code generation tasks, including reasoning (Section 6.3) and code synthesis and understanding (Section 6.4).</p>
Known Caveats	<p>Gopher (Rae et al., 2021a) describes safety benefits and safety risks associated with large language models, including PaLM. These risks include uses of language models for language generation in harmful or deceitful settings.</p> <p>PaLM should not be used for downstream applications without a prior assessment and mitigation of the safety and fairness concerns specific to the downstream application. In particular, we recommend focusing mitigation efforts at the downstream application level rather than at the pretrained level.</p>



# Model Card in the PaLM (2022) paper II

Example

System Type	
System Description	This is a standalone model.
Upstream Dependencies	None.
Downstream Dependencies	None.
Implementation Frameworks	
Hardware & Software: Training	Hardware: TPU v4 ( <a href="#">Jouppi et al., 2020</a> ). Software: T5X (t5x, 2021), JAX ( <a href="#">Bradbury et al., 2018</a> ), Pathways ( <a href="#">Barham et al., 2022</a> ). For details, see Training Infrastructure (Section 4).
Hardware & Software: Deployment	Hardware: TPU v4 ( <a href="#">Jouppi et al., 2020</a> ). Software: T5X (t5x, 2021).
Compute Requirements	Reported in Compute Usage (Section B).
Model Characteristics	
Model Initialization	The model is trained from a random initialization.
Model Status	This is a static model trained on an offline dataset.
Model Stats	The largest PaLM model has 540 billion dense parameters. We have also trained 8 billion and 62 billion parameter models.

# Model Card in the PaLM (2022) paper II

Example

System Type	
System Description	This is a standalone model.
Upstream Dependencies	None.
Downstream Dependencies	None.
Implementation Frameworks	
Hardware & Software: Training	Hardware: TPU v4 (Jouppi et al., 2020). Software: T5X (t5x, 2021), JAX (Bradbury et al., 2018), Path-

Finally, we report the net tCO<sub>2</sub>e emitted by training PaLM-540B following Patterson et al. (2021). We trained PaLM 540B in Google's Oklahoma datacenter, which has PUE of 1.08. The Oklahoma datacenter is substantially powered by wind and other renewable energy sources, and operated on 89% carbon-free energy<sup>21</sup> during the time period that the PaLM-540B was trained, with 0.079 tCO<sub>2</sub>e/MWH.<sup>22</sup> We trained PaLM-540B on 6144 TPU v4 chips for 1200 hours and 3072 TPU v4 chips for 336 hours including some downtime and repeated steps. Using 378.5W measured system power per TPU v4 chip, this leads to a total effective emissions of 271.43 tCO<sub>2</sub>e. To put this in perspective, total emissions of a direct round trip of a single passenger jet between San Francisco and New York (JFK) is estimated to be 180 tCO<sub>2</sub>e (Patterson et al., 2021), and total emissions for GPT-3 are estimated to be 552 tCO<sub>2</sub>e (Patterson et al., 2021). All of the energy use and emissions for PaLM training and the experiments described in this paper are compensated with renewable energy sources (Sustainability, 2022).

# Model Card in the PaLM (2022) paper III

Example

Data Overview	
Training Dataset	See Datasheet (Appendix D) for the description of datasets used to train PaLM.
Evaluation Dataset	We evaluate the PaLM family of models on a wide variety of tasks. Specifically, we evaluate the models on English Natural Language Processing (NLP) tasks (Section 6.1), tasks from BIG-bench (BIG-bench collaboration, 2021), reasoning tasks (Section 6.3), code completion tasks (Section 6.4), multilingual generation and question answering tasks (Section 6.6), translation tasks (Section 6.5), and bias and toxicity benchmarks (Rudinger et al., 2018; Gehman et al., 2020).
Fine-tuning Dataset	We include finetuning results on SuperGLUE (Wang et al., 2019b), tasks from GEM (Gehrmann et al., 2021), and TyDiQA (Clark et al., 2020). We also finetune on a code dataset and share results on the finetuned model on code synthesis tasks.
Evaluation Results	
Benchmark Information	<ul style="list-style-type: none"> <li>• Fewshot: English Natural Language Processing (NLP) tasks (Section 6.1), BIG-bench (Section 6.2), Reasoning (Section 6.3), Code (Section 6.4), GEM (Section 6.6), Translation (Section 6.5), Multi-lingual Question Answering (Section 6.7)</li> <li>• Finetuning: SuperGLUE (Section 6.1.1), GEM (Section 6.6), TyDiQA (Section 6.7).</li> <li>• Responsible AI: Co-occurrence, Winogender (Section 10.1.1), Real-Toxicity (Section 10.2).</li> <li>• Data contamination (Section 8)</li> </ul>
Evaluation Results	Reported in Evaluation (Section 6).

Christopher Klamm

# Model Card in the PaLM (2022) paper IV

Example

Model Usage & Limitations	
Sensitive Use	PaLM is capable of open-ended text generation. This model should not be used for any of the unacceptable language model use cases, e.g., generation of toxic speech.
Known Limitations	PaLM is designed for research. The model has not been tested in settings outside of research that can affect performance, and it should not be used for downstream applications without further analysis on factors in the proposed downstream application.
Ethical Considerations & Risks	Reported in Ethical Considerations (Section 11).

unintended use

...

# Model Card in the PaLM (2022) paper IV

Example

Model Usage & Limitations	
Sensitive Use	PaLM is capable of open-ended text generation. This model should not be used for any of the unacceptable language model use cases, e.g., generation of toxic speech.
Known Limitations	PaLM is designed for research. The model has not been tested in settings outside of research that can affect performance, and it should not be used for downstream applications without further analysis on factors in the proposed downstream application.
Ethical Considerations & Risks	Reported in Ethical Considerations (Section 11).

unintended use

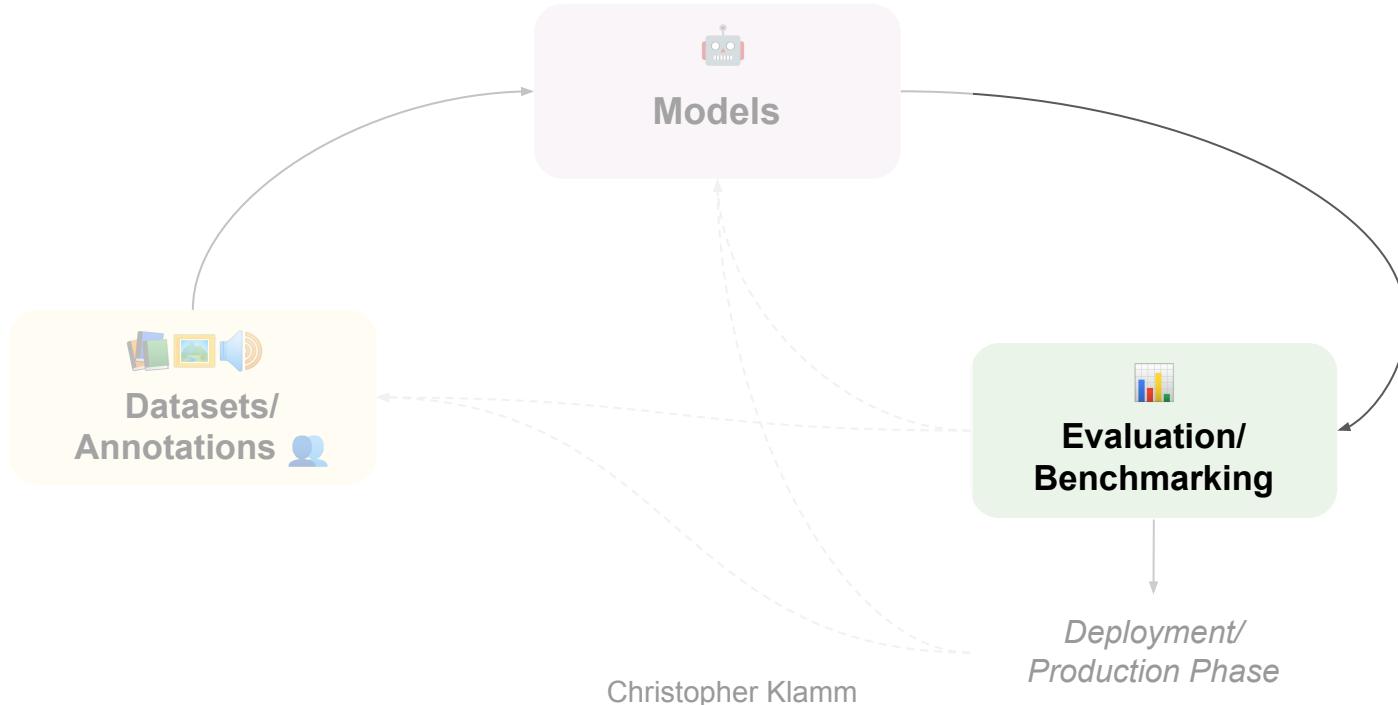


...

Our analysis reveals that our training data, and consequently PaLM, do reflect various social stereotypes and toxicity associations around identity terms. Removing these associations, however, is non-trivial; for instance, filtering off content that is deemed *toxic* by an automated tool may disproportionately exclude content about or authored by marginalized subgroups in the training data (Dodge et al., 2021). Future work should look into effectively tackling such undesirable biases in data, and their influence on model behavior. Meanwhile, any real-world use of PaLM for downstream tasks should perform further contextualized fairness evaluations to assess the potential harms and introduce appropriate mitigation and protections. ...

Chowdhery et al. 2022: 45

# Machine Learning Lifecycle/ Pipeline

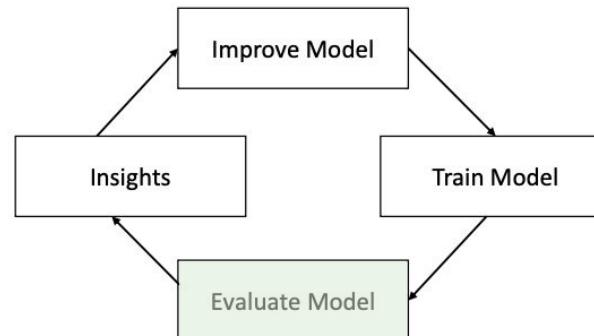


# Evaluation/ Benchmarking



# Evaluation

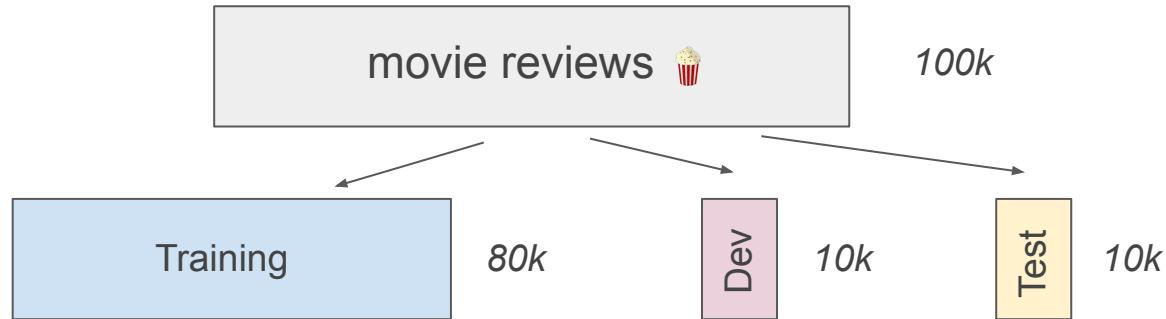
“You can’t improve what you don’t measure” (Peter Drucker)



[Reimers 2022](#)

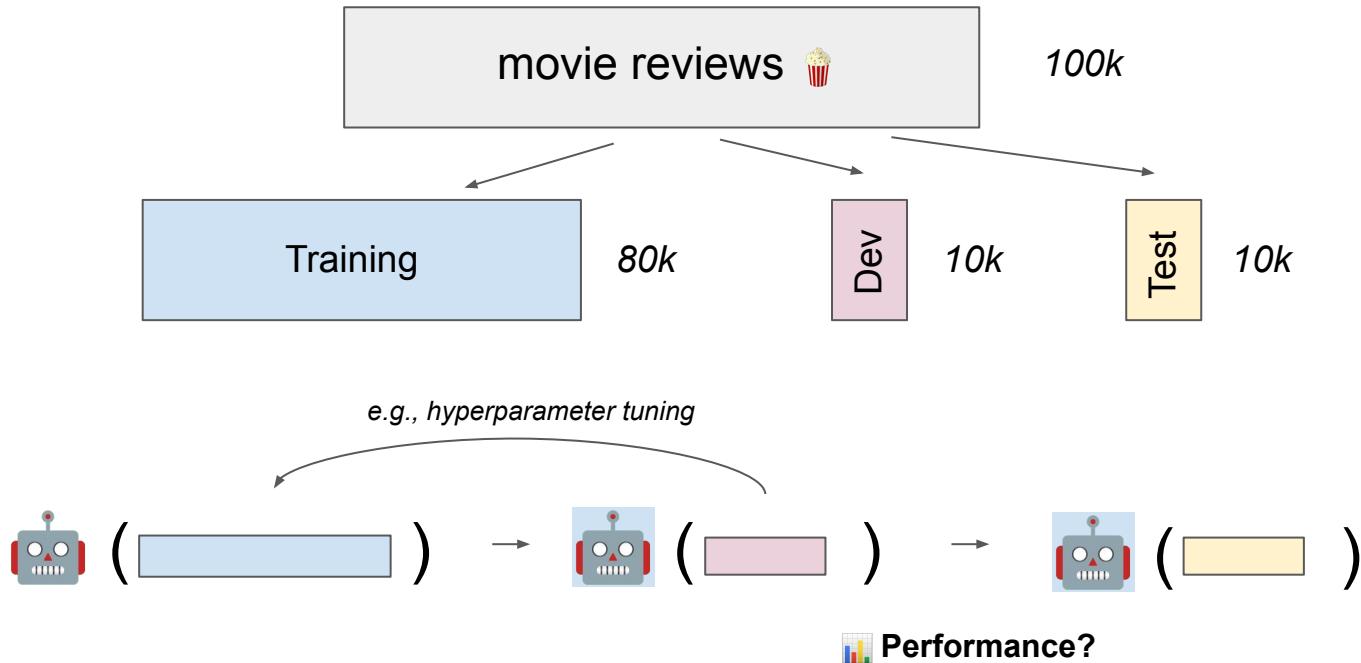


# Training/ Development/ Test Split





# Training/ Development/ Test Split





# Performance

- Accuracy

$$= \frac{TP + TN}{TP + TN + FP + FN}$$

- Precision

$$= \frac{TP}{TP + FP}$$

- Recall

$$= \frac{TP}{TP + FN}$$

$$= \frac{2 * P * R}{P + R}$$

		Predicted	
		Negative	Positive
Actual	Negative	True Negative	False Positive
	Positive	False Negative	True Positive
<i>Recall</i>		<i>Precision</i>	



# Performance

- Accuracy

$$= \frac{TP}{TP + TN + FP + FN}$$

## accuracy

Accuracy is the proportion of correct predictions among the total number of cases processed. It can be computed with:  $\text{Accuracy} = (TP + TN) / (TP + TN + FP + FN)$  Where: TP: True positive TN: True negative FP: False positive FN: False negative

- Precision

$$= \frac{TP}{TP + FP}$$

## bertscore

BERTScore leverages the pre-trained contextual embeddings from BERT and matches words in candidate and reference sentences by cosine similarity. It has been shown to correlate with human judgment on sentence-level and system-level evaluation. Moreover, BERTScore computes precision, recall, and F1 measure, which can be useful for...

- Recall

$$= \frac{TP}{TP + FN}$$

## bleu

BLEU (Bilingual Evaluation Understudy) is an algorithm for evaluating the quality of text which has been machine-translated from one natural language to another. Quality is considered to be the correspondence between a machine's output and that of a human: "the closer a machine translation is to a professional human translation, the better it is"...

- F1

$$= \frac{2 * P * R}{P + R}$$

## brier\_score

The Brier score is a measure of the error between two probability distributions.

## cer

Character error rate (CER) is a common metric of the performance of an automatic speech

## Predicted

gative

Positive

Negative

False Positive

Negative

True Positive

**Precision**



# Experimental Results Checklist ([Dodge et al. 2019](#))



“[...] test-set performance scores alone are insufficient for drawing accurate conclusions about which model performs best.” ([Dodge et al. 2019: 1](#))

✓ **For all reported experimental results**

- Description of computing infrastructure
- Average runtime for each approach
- Details of train/validation/test splits
- Corresponding validation performance for each reported test result
- A link to implemented code

✓ **For experiments with hyperparameter search**

- Bounds for each hyperparameter
- Hyperparameter configurations for best-performing models
- Number of hyperparameter search trials
- The method of choosing hyperparameter values (e.g., uniform sampling, manual tuning, etc.) and the criterion used to select among them (e.g., accuracy)
- Expected validation performance, as introduced in §3.1, or another measure of the mean and variance as a function of the number of hyperparameter trials.



Dodge et al.

2019		
Computing infrastructure	GeForce GTX 1080 GPU	
Number of search trials	50	
Search strategy	uniform sampling	
Best validation accuracy	40.5	
Training duration	39 sec	
Model implementation	<a href="http://github.com/allenai/show-your-work">http://github.com/allenai/show-your-work</a>	
Hyperparameter	Search space	Best assignment
number of epochs	50	50
patience	10	10
batch size	64	64
embedding	GloVe (50 dim)	GloVe (50 dim)
encoder	ConvNet	ConvNet
max filter size	<i>uniform-integer</i> [3, 6]	4
number of filters	<i>uniform-integer</i> [64, 512]	332
dropout	<i>uniform-float</i> [0, 0.5]	0.4
learning rate scheduler	reduce on plateau	reduce on plateau
learning rate scheduler patience	2 epochs	2 epochs
learning rate scheduler reduction factor	0.5	0.5
learning rate optimizer	Adam	Adam
learning rate	<i>loguniform-float</i> [1e-6, 1e-1]	0.0008

Table 2: SST (fine-grained) CNN classifier search space and best assignments.

[Dodge et al.  
2019](#)



# Benchmarks

“Progress in NLP has traditionally been measured through a **selection of task-level datasets** that gradually became accepted benchmarks [...]” ([Kiela et al. 2021: 2](#))

→ new **dynamic benchmarks** (as “living entities”), e.g., [Dynabench \(Kiela et al. 2021\)](#)

⚠ “On multiple benchmark setups [...], we show that the relative **performance of algorithms may be altered significantly simply by choosing different benchmark tasks**, highlighting the fragility of the current paradigms [...]” ([Dehghani et al. 2021](#))

⚠ “[...] ‘progress’ are **largely defined by performance on datasets**” but they often share “representational concerns”, “annotation artifacts”, under-specified data selection, reuse for unintended purpose, lack of standards etc. ([Paullada et al. 2020](#))

→ a need for **more comprehensive evaluations** e.g., carbon footprint, model size, fairness, robustness, etc. (see “A Critique of NLP Leaderboards”, [Ethayarajh/ Jurafsky 2020](#))



# Benchmark Checklist [\(Reimers 2022\)](#)

- **What is the intended use-case?** predicting a label, ranking of results, ...
- **Costs of Errors**
  - Research treats all errors often with equal costs
  - In production this is seldom the case
- **Human upper bound**
  - How good are humans in this task?
  - When creating a new dataset: Spend many cycles to improve human agreement
- **What else is important?**
  - Inference speed
  - Robustness
- **A benchmark must evolve**
  - As models evolve, our benchmarks must evolve!
  - Stop using outdated benchmarks
- **Restrict number of submissions**
  - The more experiments we run on a benchmark, the less likely we can trust the numbers
  - Only allow evaluation on test set in very rare cases!
  - Have a dev dataset for model development
  - If possible: use an "out-of-domain" test dataset
- **Temporal split:** Test data should be the most recent, train data the oldest
- **Diversity:** Don't test only on one task / domain etc.
- **Look for biases:**
  - What biases does your dataset have?
  - What biases does your benchmark has? e.g. only sentence pair comparison tasks

[Reimers 2022: 25-27](#)



# Benchmark Checklist ([Dehghani et al. 2021](#) based on [Gebru et al. 2021](#))

## Benchmarking checklist for reviewers and area chairs

- If there is written dissatisfaction about the author's choice of baselines, tasks, or benchmarks in the reviews, are there rationals beyond the fact that these requested datasets are "must-have" benchmarks?
- Are the reviews considering potential benefits like efficiency, fairness, and simplicity of the proposed model outside the commonly evaluated performance metrics (e.g., accuracy)?
- Are there any negative points in the reviews due to the paper proposing a method that deviates from the current trend/hype. If so, are there rational justifications for this?
- If the reviews penalizing the paper due to the proposed method not performing well only on a subset of tasks, is there enough logical elaboration on such criticism in the reviews?
- Are the reviews assessing the evaluation strategy in terms of studying the effect of different sources of variance (e.g., multiple splits, multiple random seeds, etc.)?

- If there are analyses on statistical significance testing, are they appreciated in the reviews? If there is no such analysis, are there recommendations on this provided in the reviews?
- If the paper is claiming SOTA or improvements over baselines on a benchmark, are there ablations on how much such improvement is secured by the tricks that are not tied to the main contributions?
- If the reviews are asking for more experiments, analysis, or evaluation on more benchmarks, are the potential blockers considered for such requests? E.g. those experiments being out of reach in terms of computing budget (pre-training or extremely large datasets).
- If the paper is proposing a new idea while deviating from the common paradigms, is the "out of the hype" thinking valued in the reviews as opposed to solely recognizing SOTA performance?

[Dehghani et al. 2021: 34](#)

# Responsibility

# ⚡ No free lunch

- growing **need of resources** with a bigger carbon footprint 🌐
- represent a **hegemonic worldview** 🌎 due to the used training data
- LMs are using **abusive language** 😡
- **bad actors** 😈 who abuse the possibilities of LMs
- **uncritical use** 😵 of the output (e.g., Machine Translation)
- providing **dangerous knowledge** 🚨 (e.g. tax avoidance)
- can include **personally identifiable information** 🔎
- ...



# Costs vs. 🚀 Performance?

\$2.40 vs. \$60,000

Model	Costs encoding 1M docs	Performance
<b>OpenAI GPT-3 Embedding Models</b>		
text-similarity-ada	\$800	61.86
text-similarity-babbage	\$1,200	62.62
text-similarity-curie	\$6,000	62.39
text-similarity-davinci	\$60,000	58.11
<b>Google Embedding Models</b>		
<a href="#">st5-base-1</a>	\$0.70	67.84
<a href="#">st5-large-1</a>	\$2.40	68.74
<a href="#">st5-3b-1</a>	\$6.80	69.23
<a href="#">universal-sentence-encoder-large-5</a>	\$0.35	64.51
<b>Sentence-Transformers Model</b>		
<a href="#">all-MiniLM-L6-v2</a>	\$0.12	68.06
<a href="#">all-MiniLM-L12-v1</a>	\$0.25	68.83
<a href="#">all-mpnet-base-v1</a>	\$0.70	69.98
<a href="#">all-roberta-large-v1</a>	\$2.40	70.23

[Reimers 2022](#)

\$3 vs. \$1,260,000

Model	Performance on 11 IR datasets (nDCG@10)	Cost to encode 6M Wikipedia articles
<b>OpenAI Models</b>		
cpt-text S (Ada)	49.0	\$17,000
cpt-text M (Babbage)	50.5	\$25,000
cpt-text L (Curie)	50.9	\$126,000
cpt-text XL (Davinci)	52.8	\$1,260,000
<b>Freely available models</b>		
SpladeV2	52.7	\$3

+ more memory and time is needed

[Reimers 2022](#)

→ e.g., [Zhang et al. 2022](#) “We show that OPT-175B is comparable to GPT-3, while requiring only 1/7th the carbon footprint to develop.” p. 1

Consumption	CO <sub>2</sub> e (lbs)
Air travel, 1 passenger, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000

**Training one model (GPU)**

NLP pipeline (parsing, SRL)	39
w/ tuning & experimentation	78,468
Transformer (big)	192
w/ neural architecture search	626,155

What is the carbon footprint of streaming Netflix? 

Model	Hardware	Power (W)	Hours	kWh·PUE	CO <sub>2</sub> e	Cloud compute cost
Transformer <sub>base</sub>	P100x8	1415.78	12	27	26	\$41–\$140
Transformer <sub>big</sub>	P100x8	1515.43	84	201	192	\$289–\$981
ELMo	P100x3	517.66	336	275	262	\$433–\$1472
BERT <sub>base</sub>	V100x64	12,041.51	79	1507	1438	\$3751–\$12,571
BERT <sub>base</sub>	TPUv2x16	—	96	—	—	\$2074–\$6912
NAS	P100x8	1515.43	274,120	656,347	626,155	\$942,973–\$3,201,722
NAS	TPUv2x1	—	32,623	—	—	\$44,055–\$146,848
GPT-2	TPUv3x32	—	168	—	—	\$12,902–\$43,008

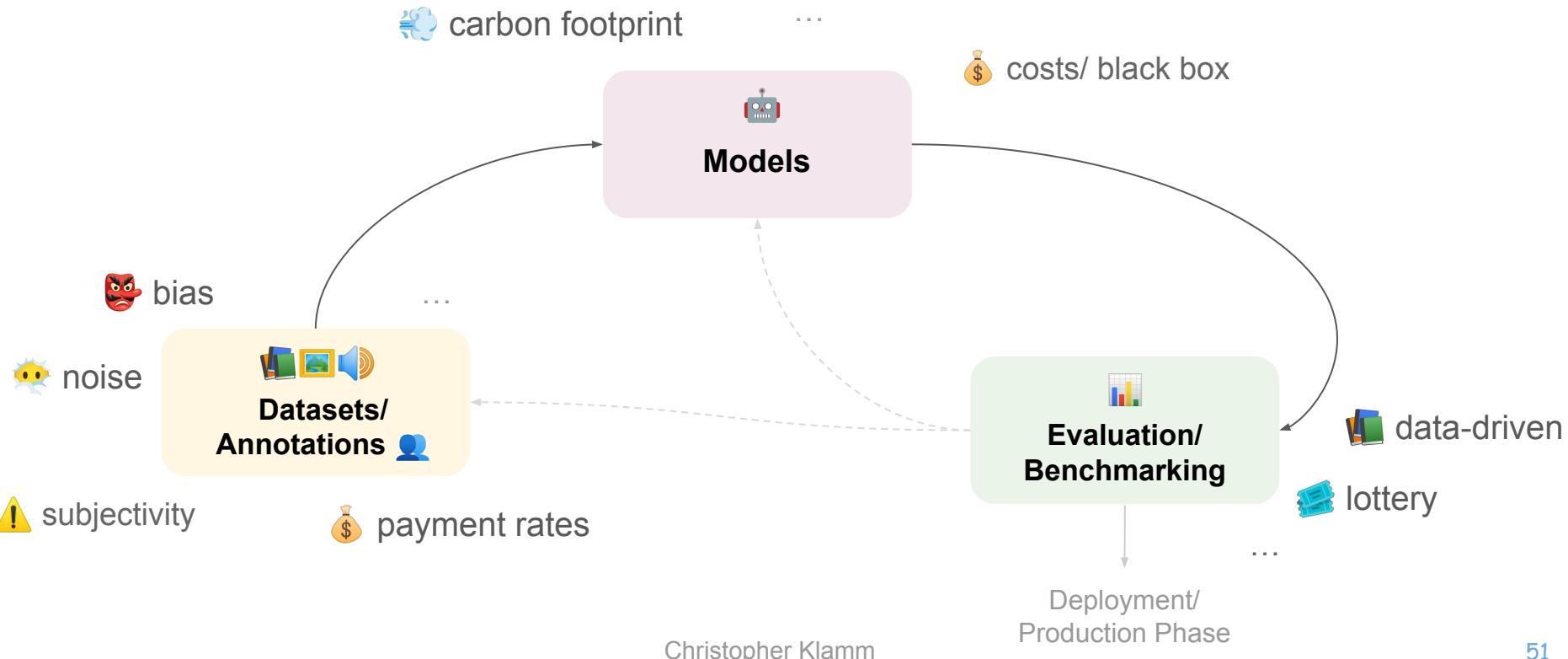
Year	Model	# of Parameters	Dataset Size
2019	BERT [39]	3.4E+08	16GB
2019	DistilBERT [113]	6.60E+07	16GB
2019	ALBERT [70]	2.23E+08	16GB
2019	XLNet (Large) [150]	3.40E+08	126GB
2020	ERNIE-GEN (Large) [145]	3.40E+08	16GB
2019	RoBERTa (Large) [74]	3.55E+08	161GB
2019	MegatronLM [122]	8.30E+09	174GB
2020	T5-11B [107]	1.10E+10	745GB
2020	T-NLG [112]	1.70E+10	174GB
2020	GPT-3 [25]	1.75E+11	570GB
2020	GShard [73]	6.00E+11	—
2021	Switch-C [43]	1.57E+12	745GB



track your models carbon footprint → <https://codcarbon.io/index.html>

big, bigger,  
biggest 

# ⚠ Machine Learning Lifecycle/ Pipeline



# Democratize Machine Learning w/ HuggingFace 😊 ?

# HuggingFace

“We're on a journey to advance and **democratize artificial intelligence** through open source and open science.”  
[\(HuggingFace 2021\)](#)

- NLP startup
- API
- open-source community 
- **models**  >80k,
- **datasets**  >12k,
- **metrics** 
- ... and many more!

## Philosophy

The acceleration in Artificial Intelligence (AI) and Natural Language Processing (NLP) will have a **fundamental impact on society**, as these technologies are at the core of the tools we use on a daily basis. A considerable part of this effort currently stems in NLP from training increasingly larger language models on increasingly larger quantities of texts.

Unfortunately, the resources necessary to create the best-performing models are found mainly in the hands of big technology giants. The stranglehold on this transformative technology poses some problems, from a research advancement, environmental, ethical and societal perspective. . . .

[BigScience Project](#) (2021-2022)

- <https://huggingface.co>
- <https://huggingface.co/course/> 



## 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 + 940

## Languages

- en
- es
- fr
- de
- zh
- sv
- ru
- fi + 173

## Licenses

- apache-2.0
- mit
- cc-by-4.0 + 33

## Other

- AutoTrain Compatible
- Eval Results
- Carbon Emissions

## Models 40,912

Search Models

↑ Sort: Most Downloads

## gpt2

Text Generation · Updated May 19, 2021 · ↓ 59.9M · ⚡ 84

## 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

## roberta-large

Fill-Mask · Updated May 21, 2021 · ↓ 2.71M · ⚡ 30

## sentence-transformers/paraphrase-MiniLM-L6-v2

Sentence Similarity · Updated Aug 30, 2021 · ↓ 2.03M · ⚡ 10

## roberta-large-mnli

Text Classification · Updated May 20, 2021 · ↓ 1.77M · ⚡ 17

## distilgpt2

Text Generation · Updated May 21, 2021 · ↓ 28.7M · ⚡ 48

## roberta-base

Fill-Mask · Updated Jul 6, 2021 · ↓ 12M · ⚡ 25

## distilbert-base-uncased

Fill-Mask · Updated Aug 29, 2021 · ↓ 4.82M · ⚡ 56

## xlm-roberta-large-finetuned-conll03-english

Token Classification · Updated Oct 12, 2020 · ↓ 4.27M · ⚡ 14

## bert-base-chinese

Fill-Mask · Updated May 18, 2021 · ↓ 3.41M · ⚡ 84

## bert-base-cased

Fill-Mask · Updated Sep 6, 2021 · ↓ 2.83M · ⚡ 17

## xlm-roberta-base

Fill-Mask · Updated Mar 4 · ↓ 2.56M · ⚡ 25

## bert-base-multilingual-cased

Fill-Mask · Updated May 18, 2021 · ↓ 1.95M · ⚡ 26

## deepset/roberta-base-squad2

Question Answering · Updated Feb 24 · ↓ 1.62M · ⚡ 58

# Model Card for FLAN-T5 on HuggingFace

**Model Card for FLAN-T5 XXL**

**Instruction finetuning**

Please answer the following question.  
What is the boiling point of Nitrogen?

**Chain-of-thought finetuning**

Answer the following question by reasoning step-by-step.  
The cafeteria had 23 apples. If they used 20 for lunch and bought 6 more, how many apples do they have?

**Language model**

-320.4F

The cafeteria had 23 apples originally. They used 20 to make lunch. So they have  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ .

**Multi-task instruction finetuning (1.8K tasks)**

**Inference: generalization to unseen tasks**

Q: Can Geoffrey Hinton have a conversation with George Washington? Give the rationale before answering.

Geoffrey Hinton is a British-Canadian computer scientist born in 1947. George Washington died in 1799. Thus, they could not have had a conversation together. So the answer is "no".

**Table of Contents**

0. [TL;DR](#)
1. [Model Details](#)
2. [Usage](#)
3. [Uses](#)
4. [Bias, Risks, and Limitations](#)

**Downloads last month**  
1,499

**Hosted inference API**

**Text2Text Generation**

Answer the following yes/no question by reasoning step-by-step. Can you write a whole Haiku in a single tweet?

**Compute** **X+Enter** 0,4

This model can be loaded on the Inference API on-demand.

**JSON Output** **Maximize**

**Datasets used to train google/flan-t5-xxl**

- esnli
- lambada
- taskmaster2

22.10.2022

Christopher Klamm

# Model Card for FLAN-T5 on HuggingFace



## Bias, Risks, and Limitations

The information below in this section are copied from the model's [official model card](#):

*"Language models, including Flan-T5, can potentially be used for language generation in a harmful way, according to Rae et al. (2021). Flan-T5 should not be used directly in any application without a prior assessment of safety and fairness concerns specific to the application."*

**"... should not be used directly in any application ..."**

## Ethical considerations and risks

*"Flan-T5 is fine-tuned on a large corpus of text data that was not filtered for explicit content or assessed for existing biases. As a result the model itself is potentially vulnerable to generating equivalently inappropriate content or replicating inherent biases in the underlying data."*

## Known Limitations

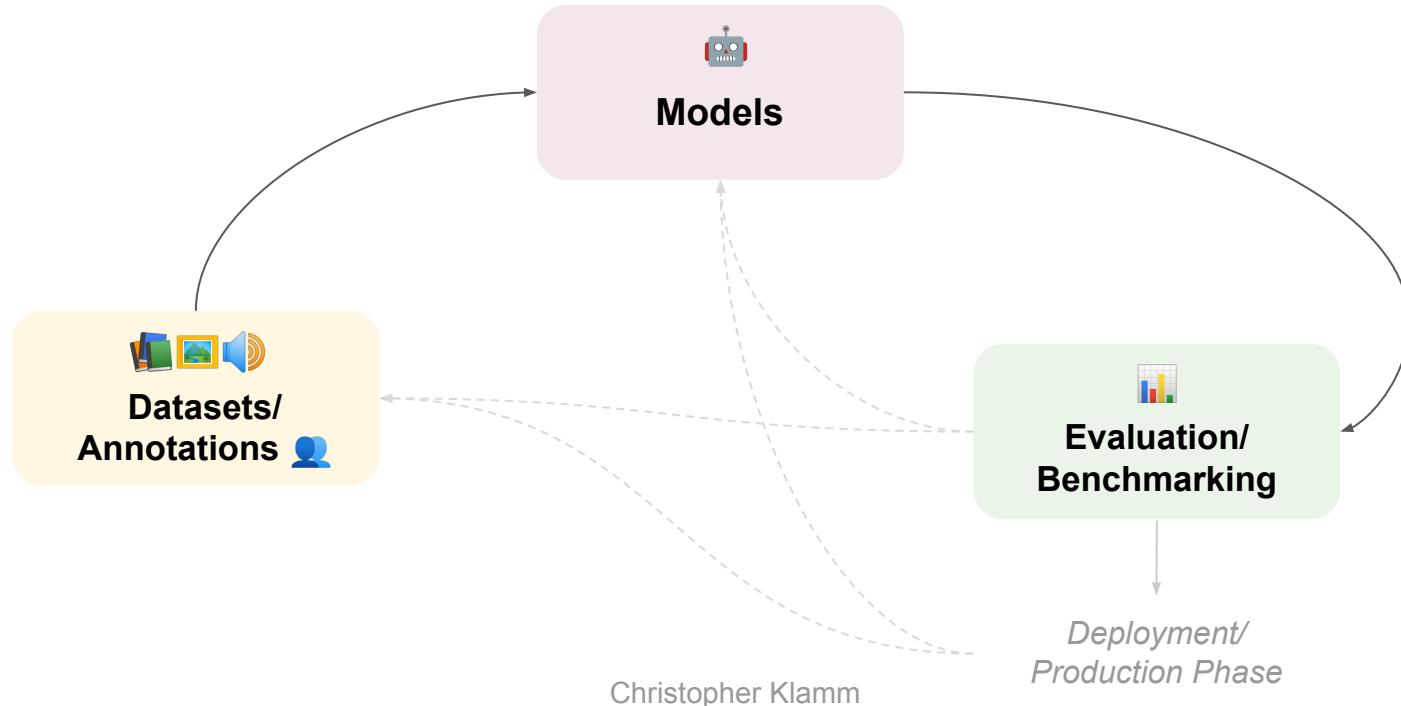
*"Flan-T5 has not been tested in real world applications."*

## Sensitive Use:

*"Flan-T5 should not be applied for any unacceptable use cases, e.g., generation of abusive speech."*

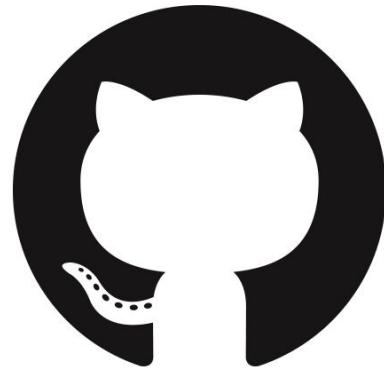
22.10.2022

# Machine Learning Lifecycle/ Pipeline



# Questions

Christopher Klam  
@chkla



[github.com/chkla/NLP-Pipeline-Talk](https://github.com/chkla/NLP-Pipeline-Talk)