


UNIVERSITAS
ATMA JAYA YOGYAKARTA
serviens in lumine veritatis




How to Do AI Research?


Technology and Innovation for Culture, Education, Language, and Society
Research Group
Y. Sigit Purnomo W.P., Ph.D.

Onsite Training Riset Berbasis Artificial Intelligence
27 Agustus 2025

1



UNIVERSITAS
ATMA JAYA YOGYAKARTA
serviens in lumine veritatis




Outline


- There Is No AI Without Data
- Everyone Wants To Do The Model Work, Not The Data Work
- AI Modelling
- Research Example

Excellence, Inclusive, Humanist & Integrity www.ujay.ac.id

2



UNIVERSITAS
ATMA JAYA YOGYAKARTA
serviens in lumine veritatis



There Is No AI Without Data

3

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

There Is No AI Without Data

- In general, it is nothing new that data preparation and data quality are key for AI and data analytics, as there is no AI without data
- There Is No AI Without Data: Industry experiences on the data challenges of AI and the call for a data ecosystem for industrial enterprises
- Christoph Gröger. 2021. **There is no AI without data**. Commun. ACM 64, 11 (November 2021), 98–108. <https://doi.org/10.1145/3448247>

Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

4

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

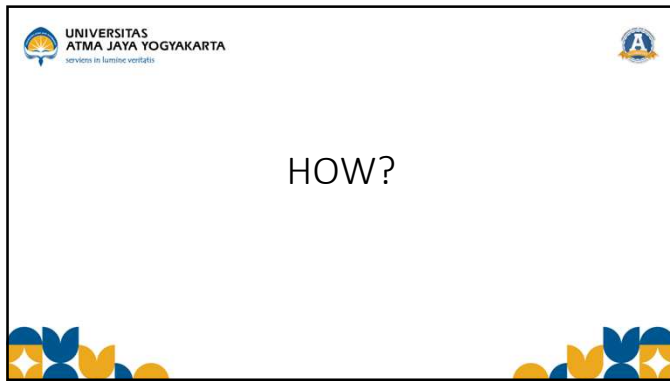
SO WHAT?

5

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

FIND THE DATA!

6



7



8

UNIVERSITAS
ATMA JAYA YOGYAKARTA
serviens in lumine veritatis

Data Source

- Culture
 - Ancient transcripts
 - Inscriptions
 - Various types of fabrics and motifs
 - Cultural data related to manners
 - Audio data of regional musical instruments
 - Cultural heritage data

Excellence, Inclusive, Humanist & Integrity

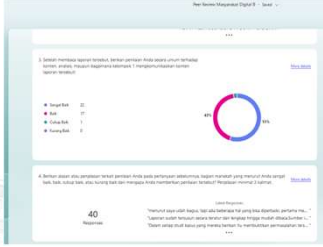
www.uajy.ac.id

9

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumine veritas

Data Source

- Education
 - Student exam/assignment results data
 - Peer review data
 - Teaching evaluation data
 - Student research (TA) results data
 - Student design/code data
 - LMS site log data
 - Labs data



Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

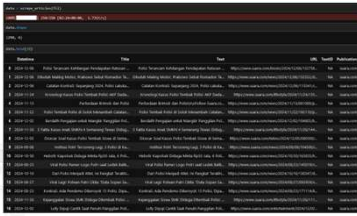
10

10

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumine veritas

Data Source

- Language
 - Text in regional languages
 - Audio/Video in regional languages
 - Text data from various news sites
 - Text data from various social media
 - Opensource application UI/UX data



Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

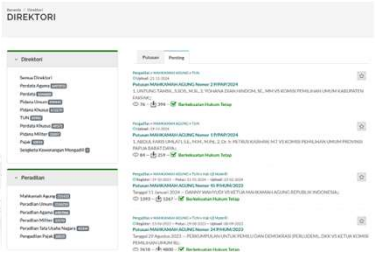
11

11

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumine veritas

Data Source

- Society
 - Data on rules or policies
 - Data related to disasters
 - Data related to climate change
 - Poverty data



Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

12

12

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumine veritas

Data Source

- **Banking & Finance**
 - financial and fiscal documents, risk assessment, and market analysis
- **Media & Content Creation**
 - multi-media data, archives and user-generated content.

Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

13

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumine veritas

Data Source

- **Conversation & Insights**
 - conversational or dialog applications and user interactions and intents
- **Healthcare & Biomedical**
 - medical documents, clinical and patient notes, pharmaceutical data and scientific research

Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

14

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumine veritas

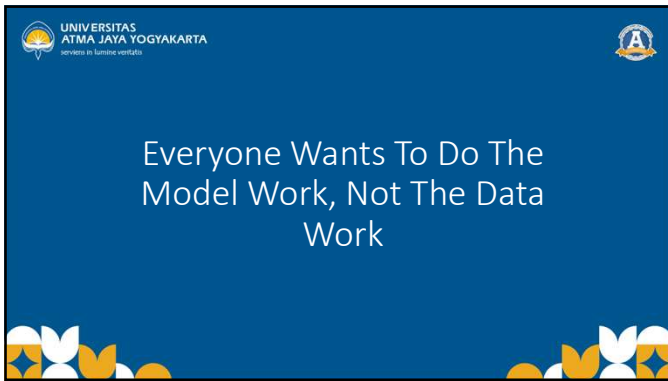
Data Source

- **Legal & Insurance**
 - legal and regulatory documents, legislative text, compliance process or domain-specific contracts.

Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

15



16

 Slide 17 contains the Universitas Atma Jaya Yogyakarta logo in the top left. The title 'Everyone wants to do the model work, not the data work' is centered. Below it, a blue bullet point states: 'Paradoxically, data is the most under-valued and de-glamorised aspect of AI'. A second bullet point reads: '“Everyone wants to do the model work, not the data work”: Data Cascades in High-Stakes AI'. A third bullet point lists authors and a citation: 'Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M Arayo. 2021. “Everyone wants to do the model work, not the data work”: Data Cascades in High-Stakes AI. In Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (CHI ’21). Association for Computing Machinery, New York, NY, USA, Article 39, 1–15. <https://doi.org/10.1145/3411764.3445518>'. The bottom of the slide features the text 'Excellence, Inclusive, Humanist & Integrity' and the website 'www.ujay.ac.id' next to a small yellow square with the number '17'.

17



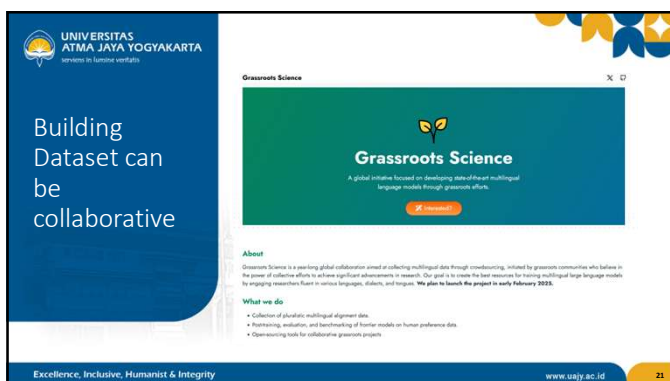
18



19



20



21

UNIVERSITAS ATMA JAYA YOGYAKARTA
servans in lumine veritas

Building Dataset can be collaborative

What's Next?

SEA-VL: Developing Culturally Relevant Vision-Language Models for Southeast Asia

Following the success of our SEA-Crowd project, we're excited to announce SEA-VL, a new open source initiative to create high quality vision-language datasets specifically for Southeast Asian (SEA) languages! We're calling on contributors to help us build a SEA-specific vision-language model.

SEA-VL
Developing Culturally Relevant Vision-Language Models for Southeast Asia

SEA-VL is a big initiative, so we have decided to split it into two phases. In **Phase 1** of SEA-VL, we're looking for self-salors, culturally-relevant images with descriptions about the shared image. This will be cleaned and compiled into a comprehensive open-access SEA-relevant image dataset. This dataset will serve as the foundation for Phase 2, where we'll develop instruction-tuning VL datasets and build a SEA-specific vision language model (VLM) using the constructed dataset.

Phase 1 is open from 11 Nov 2024 to 15 Feb 2025.

Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

22

22

UNIVERSITAS ATMA JAYA YOGYAKARTA
servans in lumine veritas

Building Dataset can be collaborative

Projects

SEA-Crowd, SEA-VL, SEA-Chat

Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

23

23

UNIVERSITAS ATMA JAYA YOGYAKARTA
servans in lumine veritas

Building Dataset can be collaborative

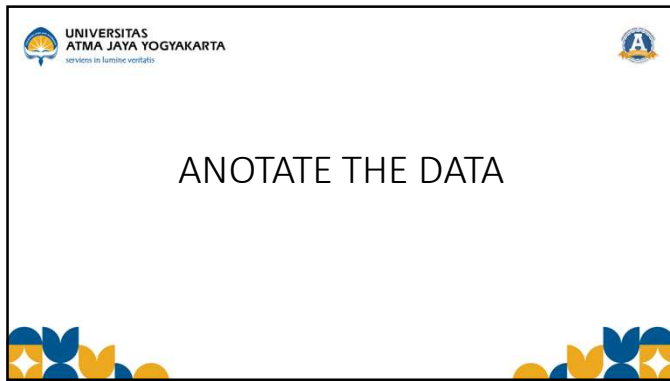
NOW, WE ALREADY HAVE THE DATA, THEN?

Excellence, Inclusive, Humanist & Integrity

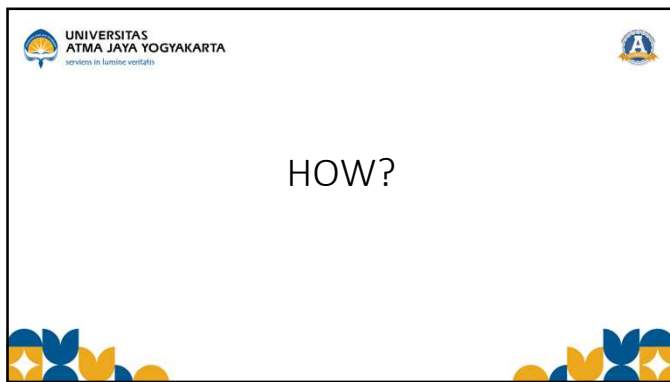
www.ujay.ac.id

24

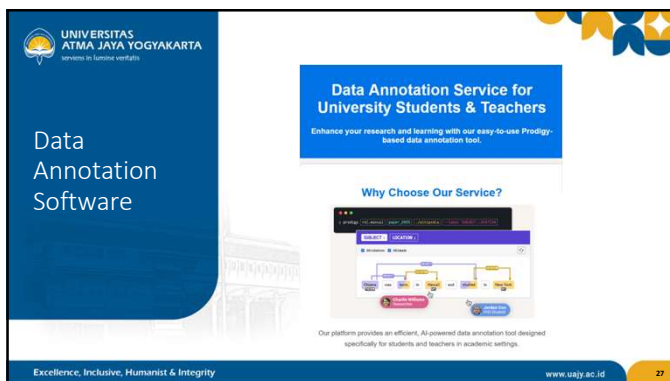
24



25



26

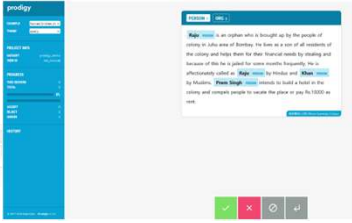


27

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

Prodigy Features

- **Information Extraction**
 - Extracting structured data like names, key phrases or relations is a crucial task for many business applications.



Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

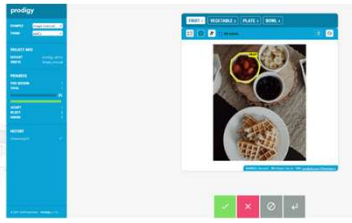
28

28

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

Prodigy Features

- **Computer Vision**
 - Prodigy makes it easy to bring structure to your images and build custom AI systems for classification, segmentation and object detection.



Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

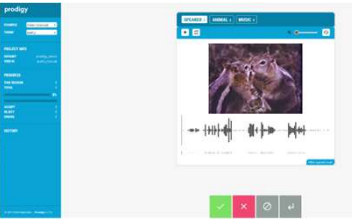
29

29

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

Prodigy Features

- **Audio & Video Annotation**
 - Prodigy makes it easy to bring structure to your audio and video data, and build custom AI systems for classification, segmentation and transcription.



Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id


30

30

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

Prodigy Features

- **Prompt Engineering**
 - Prodigy's prompt engineering workflows help you find and evaluate the best Large Language Model prompts for any custom use case.



Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id


31

31

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

Prodigy Features

- **Language Model Training**
 - Distill knowledge from large models into smaller, faster and fully private pipelines for your use case that you can run cheaply and efficiently in-house.



Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id


32

32

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

Prodigy Features

- **Custom Workflow**
 - Customize Prodigy for your specific use case and implement powerful automated workflows, interfaces and integrations.



Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

33

33

[illegible]

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in human veritas

Data Text – POS Tag

Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

37

37

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in human veritas

Data Text – Dependency Parsing

Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

38

38

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in human veritas

Data Text – Relations

Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

39

39

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in future veritas

Data Text – Text Comparison

prodigy

Nach dem Einlegen der Patrone in den Pen muss die Suspension bei einer Lagerung unterhalb von +30°C innerhalb von 28 Tagen aufgeführt werden.

After insertion of the cartridge in a pen, the suspension should be used **within 28 days** when stored **below +30°C**.

✓ ✗ ↺ ↻

Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

40

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in future veritas

Data Text – Question Answer

prodigy

What are the best education portals in India?

☐ Which is a best educational portal in india?

☐ Is the Earth flat?

✓ ✗ ↺ ↻

Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

41

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in future veritas

Data Gambar – Segmentation

prodigy

VEGETABLE | PLATE | BOWL

✓ ✗ ↺ ↻

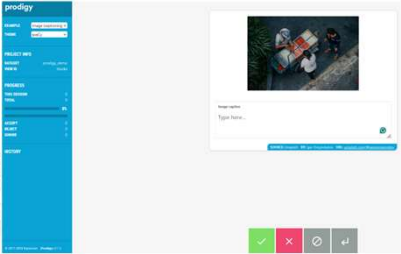
Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

42

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in future veritas

Data Gambar – Captioning



Excellence, Inclusive, Humanist & Integrity

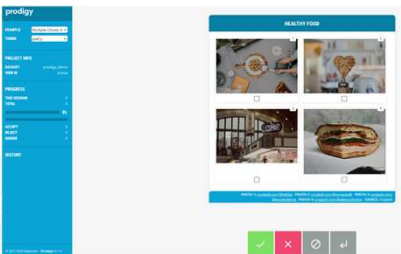
www.ujay.ac.id

43

43

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in future veritas

Data Gambar – Classification



Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

44

44

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in future veritas

Data Audio – Classification



Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

45

45



46



47



48

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

Dokumen PDF

EVENT 1 | PLACE 2

The Student StarCraft AI Tournament (SSCAIT) is the StarCraft AI competition with the highest number of total participants. It started as a part of an AI course at **Cemerlang University**, and initial seasons included several dozen student submissions from this course, in addition to submissions from across the globe. Since then, SSCAIT started accepting non-student participants and team submissions. There are three fundamental differences between SSCAIT and the remaining two competitions:

TITLE: The Current State of StarCraft AI Competitions and Botz (Cortis & Churchill, 2020); PAGE: 2

Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

49

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

Dataset Example

- Huggingface Dataset
 - <https://huggingface.co/datasets>
 - https://sbert.net/docs/sentence_transformer/dataset_overview.html
- CultureBank
 - <https://culturebank.github.io/data>

1. Description Extraction
2. Description Clustering
3. Post-processing

1.1 Cultural Belief Classifier
1.2 Cultural Belief Classifier
1.3 Cultural Belief Classifier
1.4 Cultural Belief Classifier
1.5 Cultural Belief Classifier
1.6 Cultural Belief Classifier
1.7 Cultural Belief Classifier
1.8 Cultural Belief Classifier
1.9 Cultural Belief Classifier
1.10 Cultural Belief Classifier

2.1 Clustering
2.2 Clustering
2.3 Clustering
2.4 Clustering
2.5 Clustering
2.6 Clustering
2.7 Clustering
2.8 Clustering
2.9 Clustering
2.10 Clustering

3.1 Agreement Classifier
3.2 Agreement Classifier
3.3 Agreement Classifier
3.4 Agreement Classifier
3.5 Agreement Classifier
3.6 Agreement Classifier
3.7 Agreement Classifier
3.8 Agreement Classifier
3.9 Agreement Classifier
3.10 Agreement Classifier

Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

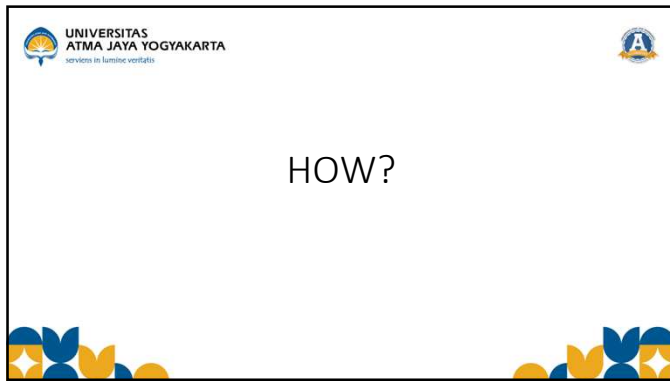
50

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

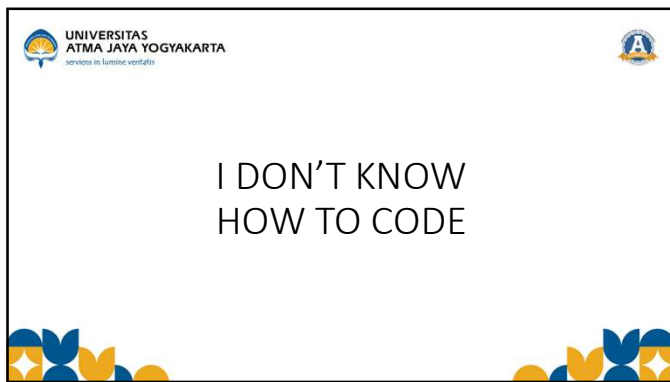
AI Modelling

Excellence, Inclusive, Humanist & Integrity

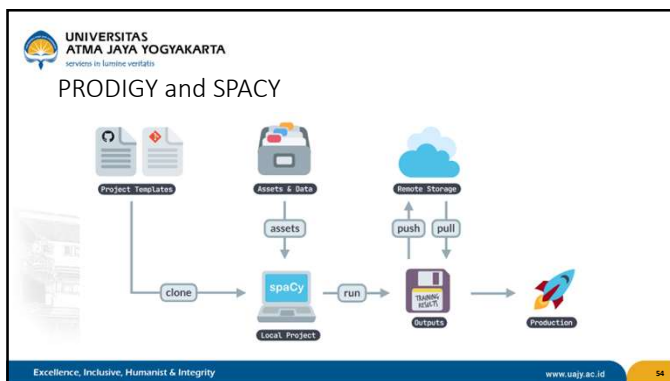
51




52



53



54




UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

PRODIGY and SPACY

- spaCy projects let you manage and share end-to-end spaCy workflows for different use cases and domains, and orchestrate training, packaging and serving your custom pipelines.
 - You can start off by cloning a pre-defined project template, adjust it to fit your needs, load in your data, train a pipeline, export it as a Python package, upload your outputs to a remote storage and share your results with your team
 - <https://github.com/explosion/projects?tab=readme-ov-file>

Excellence, Inclusive, Humanist & Integrity www.ujay.ac.id 55

55




UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

PRODIGY and SPACY

- Clone the project template you want to use
 - `python -m weasel clone tutorials/ner_fashion_brands`
- Install any project requirements.
 - `cd ner_fashion_brands`
 - `python -m pip install -r requirements.txt`
- Fetch assets (data, weights) defined in the project.yml
 - `python -m weasel assets`

Excellence, Inclusive, Humanist & Integrity www.ujay.ac.id 56

56



UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

PRODIGY and SPACY

- Run a command defined in the project.yml.
 - `python -m weasel run preprocess`
- Run a workflow of multiple steps in order.
 - `python -m weasel run all`
- Run the visualization
 - `python -m spacy project run visualize`
- Adjust the template for your specific use case, load in your own data, adjust the settings and model and share the result with your team.

Excellence, Inclusive, Humanist & Integrity www.ujay.ac.id 57

57

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

Zero/Few-shot Prompting

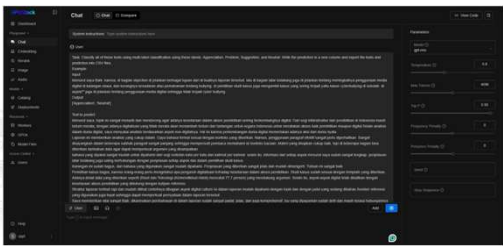
- **Zero-shot prompting** asks a large language model (LLM) to perform a task using only an instruction (no examples)
 - It relies on the model's prior, instruction-following ability
- **Few-shot prompting** adds a handful of in-context examples (input→output pairs) inside the prompt, so the model infers the task pattern on the fly: no parameter updates, just smarter conditioning

Excellence, Inclusive, Humanist & Integrity www.uajy.ac.id 58

58

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

Zero/Few-shot Prompting

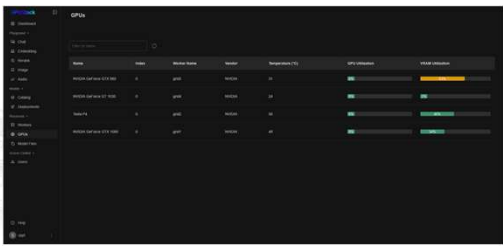


Excellence, Inclusive, Humanist & Integrity www.uajy.ac.id 59

59

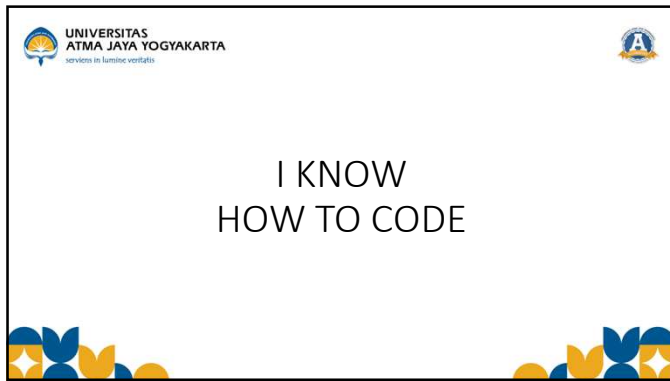
UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

Zero/Few-shot Prompting

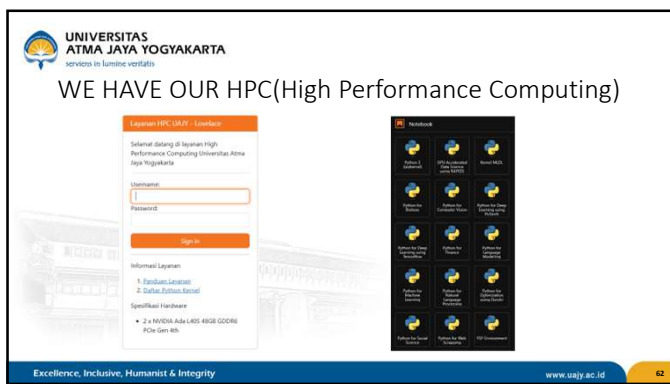


Excellence, Inclusive, Humanist & Integrity www.uajy.ac.id 60

60



61



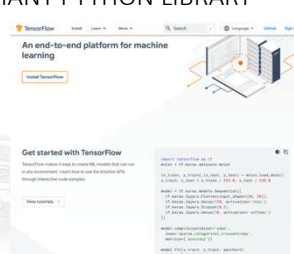
62



63

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

WE HAVE MANY PYTHON LIBRARY



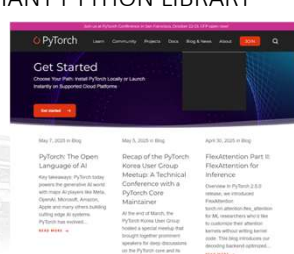
Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

64

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

WE HAVE MANY PYTHON LIBRARY



Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

65

UNIVERSITAS
ATMA JAYA YOGYAKARTA
services in lumine veritas

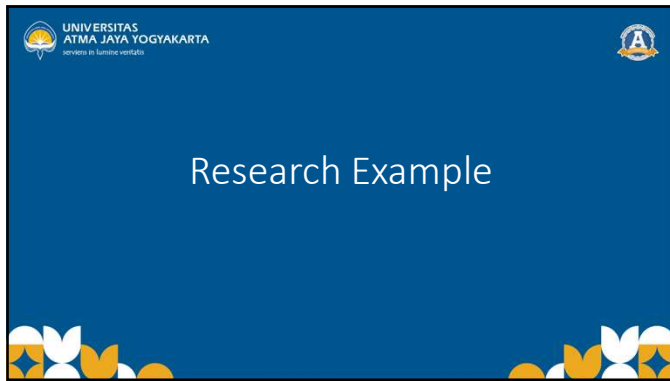
DON'T WORRY, MANY RESOURCES AVAILABLE

- Scikit Learn
 - https://scikit-learn.org/stable/auto_examples/index.html
- Keras
 - <https://keras.io/examples/>
- Tensorflow
 - <https://www.tensorflow.org/tutorials>
- PyTorch
 - <https://docs.pytorch.org/tutorials/>
- Open-Source AI Cookbook
 - <https://huggingface.co/learn/cookbook/index>

Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

66



67

UNIVERSITAS ATMA JAYA YOGYAKARTA
services to lumine veritas

Biology

- Case Study: TaxoNERD that provides deep neural network (DNN) models to recognise taxon mentions in ecological documents
 - To achieve high performance, these models usually need to be trained on a large corpus of manually annotated text. Creating such a corpus is a laborious and costly process, with the result that manually annotated corpora in the ecological domain tend to be too small to learn an accurate DNN model from scratch. To address this issue, we leverage existing models pretrained on large biomedical corpora using transfer learning.

• <https://doi.org/10.1111/2041-210X.13778>

Excellence, Inclusive, Humanist & Integrity

www.ujay.ac.id

68

UNIVERSITAS ATMA JAYA YOGYAKARTA
services to lumine veritas

Banking and Finance

- Case Study: How S&P Global is making markets more transparent with NLP
 - This case study shows how the S&P Global Commodity Insights team built and shipped impressively efficient information extraction pipelines for real-time commodities trading insights in a high-security environment, and how they were able to achieve a 10x speed-up of their data collection and annotation workflows using human-in-the-loop distillations with LLMs, spaCy and Prodigy.

• <https://explosion.ai/blog/sp-global-commodities>

The text processing workflow and spaCy pipeline

	Global Carbon Credits	Americas Crude Oil	Asia Steel Rebar
Accuracy (F-score)	0.95	0.96	0.99
Speed (words/second)	15,730	13,908	16,075
Model Size	6.14G	6.14G	6.14G
Training Examples	1,500	1,500	1,500
Evaluation Examples	211	200	345
Data Development Time	-15d	-15d	-15d

Excellence, Inclusive, Humanist & Integrity


www.ujay.ac.id

69

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumine vorlatis

Media & Content Creation

- Case Study: How the Guardian approaches quote extraction with NLP
 - To facilitate trust, human-in-the-loop workflows are widespread in media applications as stakeholders require the ability to teach and to evaluate models through human-AI interfaces.
 - <https://explosion.ai/blog/guardian>



Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

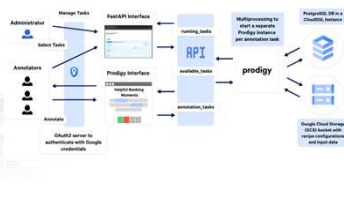
70

70

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumine vorlatis

Conversation & Insights

- Case Study: Deploying a Prodigy cloud service for Posh's financial chatbots
 - Posh focuses on developing custom NLP models trained on real-world banking conversations and custom models for each client's unique customer base and product offering.
 - <https://explosion.ai/blog/posh-prodigy-financial-chatbots>



Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

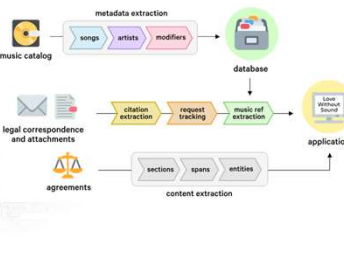
71

71

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumine vorlatis

Legal & Insurance

- Case Study: How Love Without Sound helps music industry law firms recover millions in revenue for artists
 - This case study shows how Love Without Sound built innovative AI-powered tools for the music industry and law firms specializing in royalty negotiations, and helped publishers recover hundreds of millions of dollars in lost revenue for artists using spaCy and Prodigy.
 - <https://explosion.ai/blog/love-without-sound-nlp-music-industry>



Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

72

72

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumina veritas

Computational Journalism

- Case Study: Extraction and attribution of public figures statements for journalism in Indonesia using deep learning

Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

73

73

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumina veritas

Reference

- Christoph Gröger. 2021. **There is no AI without data**. Commun. ACM 64, 11 (November 2021), 98–108. DOI: <https://doi.org/10.1145/3448247>
- Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M Aroyo. 2021. **"Everyone wants to do the model work, not the data work": Data Cascades in High-Stakes AI**. In Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (CHI '21). Association for Computing Machinery, New York, NY, USA, Article 39, 1–15. <https://doi.org/10.1145/3411764.3445518>
- <https://prodi.gy>

Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

74

74

UNIVERSITAS ATMA JAYA YOGYAKARTA
services in lumina veritas

Reference

- Yohanes Sigit Purnomo W.P., Yogan Jaya Kumar, Nur Zareen Zulkarnain, Basit Raza. 2024. **Extraction and attribution of public figures statements for journalism in Indonesia using deep learning**. Knowledge-Based System, Volume 289, 8 April 2024, 111558 - DOI: <https://doi.org/10.1016/j.knosys.2024.111558>
- Yohanes Sigit Purnomo W.P., Yogan Jaya Kumar, Nur Zareen Zulkarnain. 2022. **PFSA-ID: An Annotated Indonesian Corpus and Baseline Model of Public Figures Statements Attributions**. Global Knowledge, Memory and Communication - DOI: <https://10.1108/GKMC-04-2022-0091>

Excellence, Inclusive, Humanist & Integrity

www.uajy.ac.id

75

75



76
