# Synthetic data

Team 8: Luke Smith, Garrison Chura, Sofiya Kuzina, and Cynthia Marquez

# Problem Statement

*What is synthetic data, what is its purpose, and how effective is its use in AI modeling?*

We aim to explore the effectiveness, challenges, and applications of using synthetic data to overcome the limitations of using real-world data to train AI/ML models.

# Context

**Real data:**
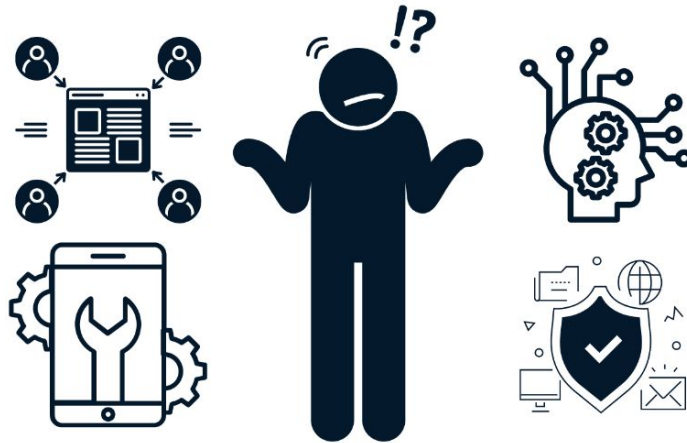
Difficult to access/laborious

Expensive

Constrained by regulations – privacy concerns

**Synthetic data:**

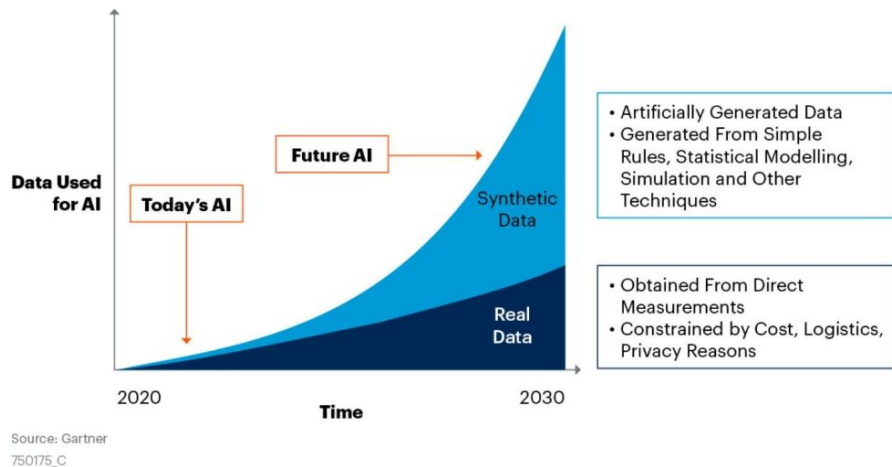Computer–generated data that is similar to real–world data

Primary purpose: increase the privacy and integrity of systems

**Source:** Creating Synthetic Data with Python
Faker Tutorial | DataCamp
(https://www.datacamp.com/tutorial/creating)

Raymond A. Mason
School of Business
WILLIAM & MARY

# Importance of Synthetic Data
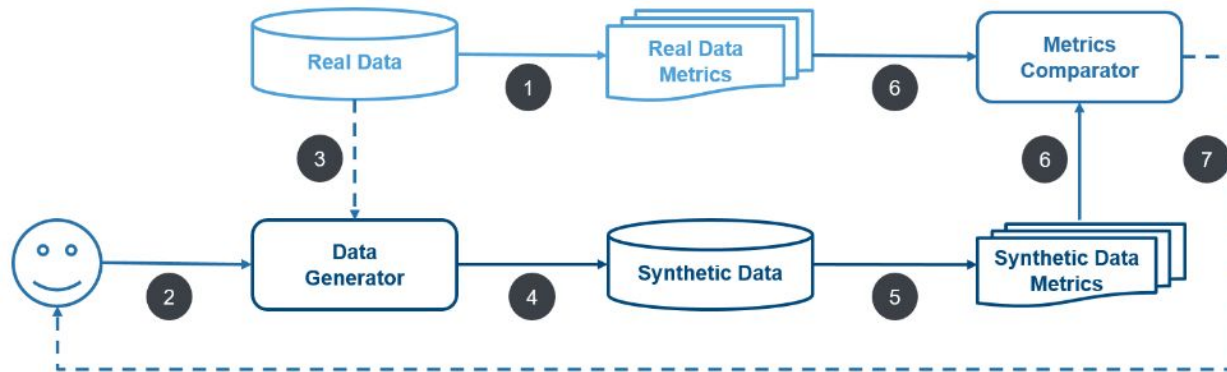
- **Privacy**

- **Cost**

- **Diversity**

- **Control**

- **Scalability**



**Source:** Creating Synthetic Data with Python Faker Tutorial | DataCamp (https://www.datacamp.com/tutorial/creating)

# Process to Generate Synthetic Financial Datasets



Step 1: Compute **metrics for the real data**

Step 2: **Develop a Generator** (may be *statistical methods* or an *agent-based simulation*)

Step 3: *(Optional)* **Calibrate the Generator** using the real data

Step 4: **Run the Generator** to generate synthetic data

Step 5: Compute **metrics for the synthetic data**

Step 6: **Compare the metrics** of the real data and synthetic data

Step 7: *(Optional)* **Refine the Generator** to improve against comparison metrics

Raymond A. Mason
School of Business
WILLIAM & MARY

# Literature Review

*"Synthetic Document Generator For Annotation-Free Layout Recognition"*
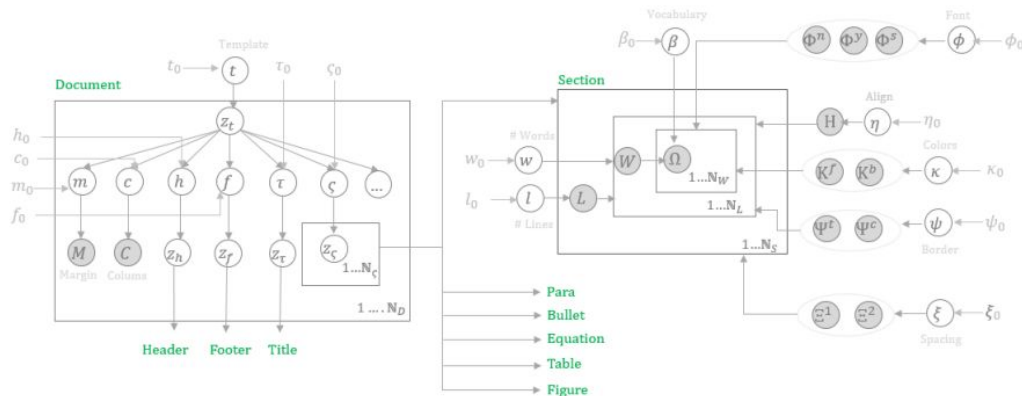
Natraj Raman, Sameena Shah and Manuela Veloso

JPMorgan AI Research Lab

# "Synthetic Document Generator for Annotation-free Layout Recognition"

- Analyze document layout

- Use Bayesian Network to make synthetic documents

- Train object detection model to predict labels for each part of a document's layout



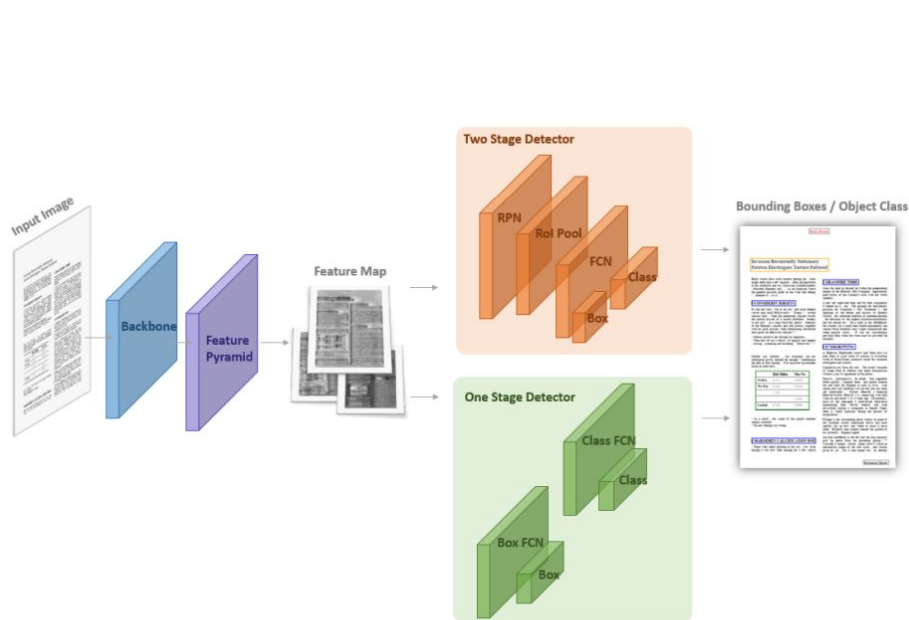**Source:** https://doi.org/10.1016/j.patcog.2022.108660

# "Synthetic Document Generator for Annotation-free Layout Recognition"

- Titles, sections, headers/footers,, tables, figures etc. help understand document content

- **Layout recognition**

  - Use object detection model

  - Input images → scaled feature maps → identify layout elements and boundaries

- **Synthetic document generator** → produces realistic documents that have labeled layout elements/spatial positions

Raymond A. Mason
School of Business
WILLIAM & MARY

# "Synthetic Document Generator for Annotation–free Layout Recognition"



**Figure 3:** Layout Recognition Model Architecture. A feature extraction network takes an image of arbitrary size as input and produces feature maps at different scales. An object detector network determines the categories and bounding boxes of the layout elements.

# "Synthetic Document Generator for Annotation–free Layout Recognition"

## *Results:*

- Train layout detectors on synthetic data → as good as real documents

- Increase number of synthetic documents → performance of real and synthetic documents converge

- "Granularity" of the layout categories could impact recognition quality

**Table 8:** Impact of training with a subset of layout categories.

| Trained Categories | Target Category | Real and Synthetic F1 Difference |
|---|---|---|
| All Categories | Section | 4.6 |
| All Categories | Table | 3.8 |
| All Categories | Figure | 2.7 |
| Only Section | Section | 5.1 |
| Section + Equation | Section | 2.6 |
| Only Table | Table | 3.6 |
| Only Figure | Figure | 1.3 |
| Table + Figure | Table | 1.7 |
| Table + Figure | Figure | 0.2 |

**Source:** https://doi.org/10.1016/j.patcog.2022.108660

Raymond A. Mason
School of Business
WILLIAM & MARY

# Code Demo

**Make_classification (sklearn.datasets)**
Function used to generate synthetic datasets for classification tasks (random n-class)
- Testing Machine Learning Algorithms
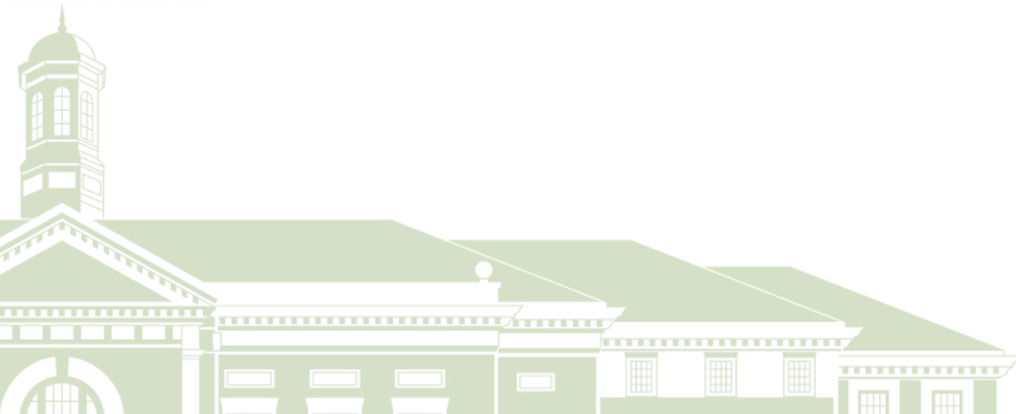- Creating Imbalanced Datasets and Data
- Synthetic Data Augmentation

# Explainability, Challenges, & Ethical Concerns

**Synthetic data can…**

- Amplify biases that exist in real-world data
- Miss complexities (such as outliers) of real data
- Not be sensitive to real-time changes
- Ethics:
  - Ownership of synthetic data from publically available data
  - Privacy → can help ensure privacy, but…
    - Data leakage – risk of individuals being identified from real data
    - Web scraping, using real data without consent

**However…**

- Synthetic datasets can as accurate than real-world data
  (Raman, S., Shah, S., & Veloso, M. (2022))

Raymond A. Mason
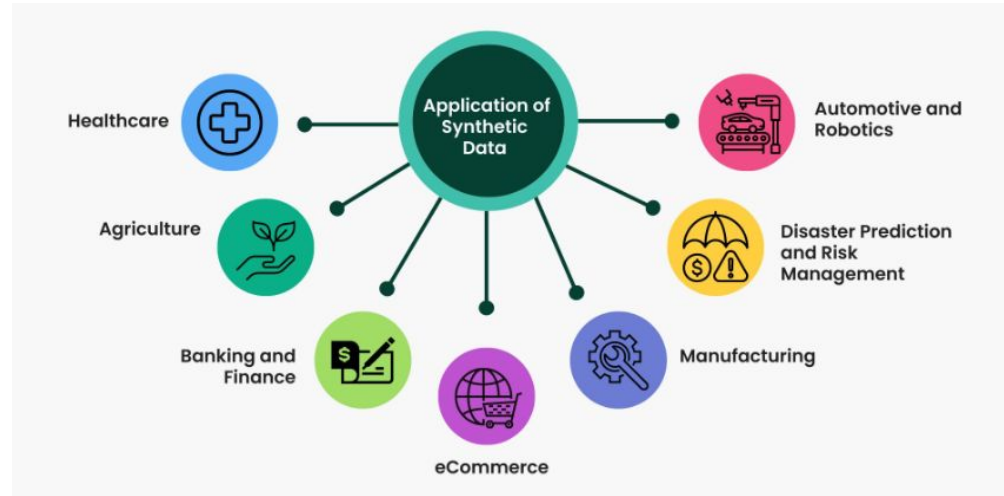School of Business
WILLIAM & MARY

# What's Next?

**Future Developments**

- Industry-specific applications
- Techniques to capture data from different sources
- New idea → spread awareness

**Concerns:**

- Address potential biases
- Ensure synthetic data represents diversity of real-world data
- Develop standard to measure synthetic data accuracy (validation)



**Source:** Synthetic Data Generation: Definition, Types, Techniques, & Tools (https://www.turing.com/kb/synthetic-data-generati)

Raymond A. Mason
School of Business
WILLIAM & MARY

# References

Awan, A. (2022, August). Creating Synthetic Data with Python Faker Tutorial. DataCamp.

https://www.datacamp.com/tutorial/creating-synthetic-data-with-python-faker-tutorial

J.P. Morgan. (n.d.). Synthetic Data. https://www.jpmorgan.com/synthetic-data

Radecic, D. (2021, January 10). How to Make Synthetic Datasets with Python: A Complete Guide for Machine

Learning. Better Data Science. https://betterdatascience.com/python-synthetic-datasets/

Raman, S., Shah, S., & Veloso, M. (2022). Synthetic document generator for annotation-free layout

recognition. Pattern Recognition, 128. https://doi.org/10.1016/j.patcog.2022.108660

Turing.com. (n.d.). Synthetic Data Generation: Definition, Types, Techniques, and Tools.

https://www.turing.com/kb/synthetic-data-generation-techniques

Zewe, A. MIT News. (2022, November 3). In machine learning, synthetic data can offer real performance

improvements. https://news.mit.edu/2022/synthetic-data-ai-improvements-1103

Raymond A. Mason
School of Business
WILLIAM & MARY