# Creating Synthetic Experts with Generative Artificial Intelligence

## Daniel M. Ringel

*dmr@unc.edu*

**Wharton Webinar:** AI and Innovation
March 1, 2024

# Multi-Purpose AI - Bigger, Better, Overkill?

## Multi-Purpose AI

→ **Applicable to many tasks**
→ **Promises to revolutionize knowledge work**

→ Unwieldy
→ Resource-hungry
→ Difficult to build and run
→ Largely closed and proprietary
→ Centralized and under the control of few

www.scientificamerican.com/article/when-it-comes-to-ai-models-bigger-isnt-always-better/

**Example:** OpenAI's GPT-4

## Specialized AI

→ **Build for individual tasks**
→ **Known for efficiency and accuracy**

→ Requires domain expertise for training

**Example:** OpenAI's Toxicity Classifier*

* used by OpenAI to maintain integrity and safety of ChatGPT

## *Why buy the whole Candy Store, when you just need a Lollypop ?*

# Specialized AI: Natural Language Classifiers

*Sort items into specific categories based on their characteristics*

**Distill intelligence from** vast amounts of **unstructured data**

- News and social media
- Customer interactions
- Reports and policies
- Internal communications

**Classifiers** can swiftly **identify constructs of interest** in data

- Specific topics
- Bias and sentiment
- Compliance
- Emotions

*Image by Midjourney*

**Versatility** of classifiers **extends** their **utility across sectors** and **functions**

# Efficacy of Classifiers relies heavily on their Training

**Training** an effective classifier typically requires *many labeled examples*

**easy** for *simple constructs*

*Straightforward; can be easily defined and measured*

- crowdsourcing (e.g., Amazon mTurk)
- fast and relatively low cost

**problematic** for more *complex constructs*

*Multifaceted with higher levels of abstraction and ambiguity*

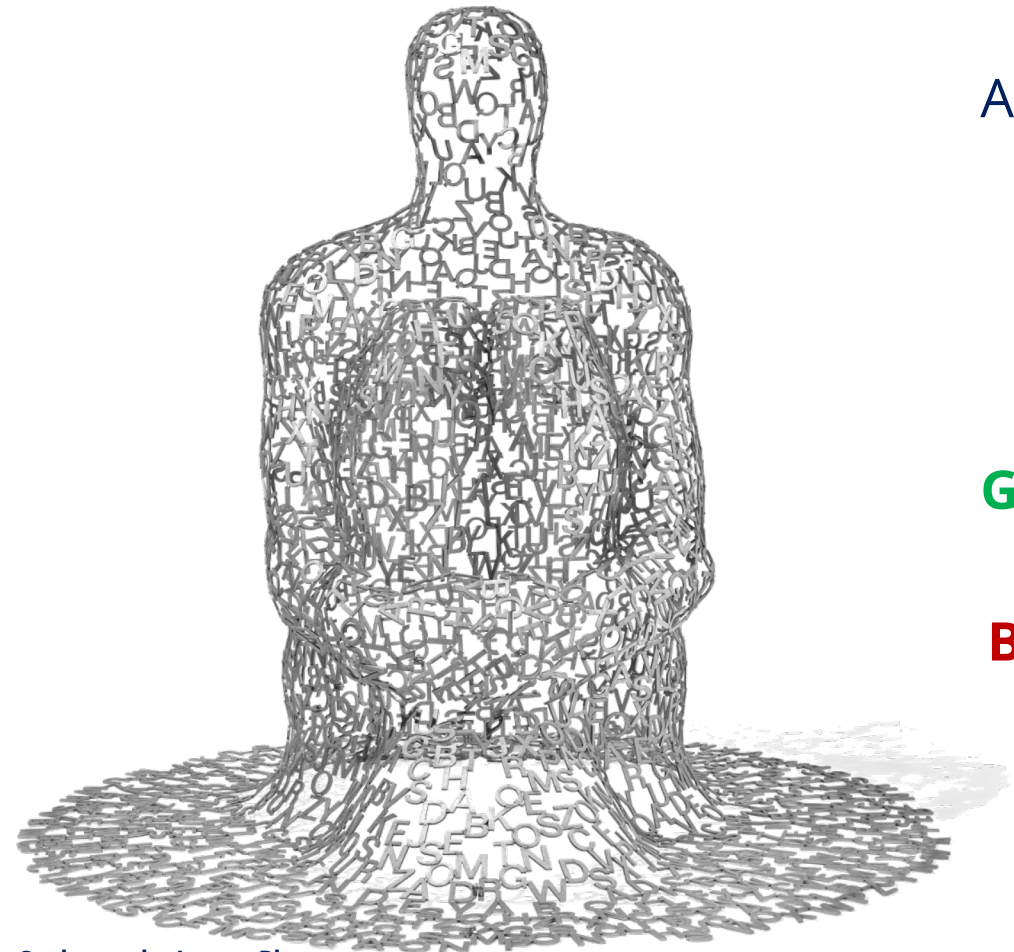- requires domain experts
- scarce and expensive resource

**Complex constructs** are often *more insightful*

*General* consumer sentiment towards a brand

**VS**

*What exactly do* consumers perceive as positive or negative about a brand
*(e.g., Marketing Mix: Product, Place, Price, Promotion)*

# Let Generative AI label Data



**Sculpture by Jaume Plensa**

Ask *generative AI* to identify complex construct of interest

- ***easily accessible*** | *e.g., OpenAI's GPT-4*
- represents ***vast body of knowledge***
- many ***theoretically founded constructs*** in training data

**Good News:** Works well!

**Bad News:** Limitations

- largely ***proprietary***
- ***slow*** and ***costly***
- ***reproducibility***
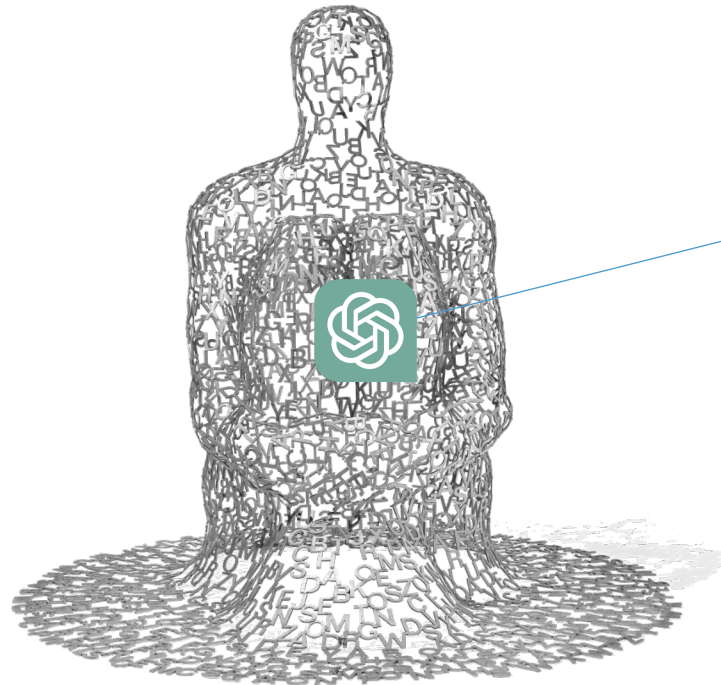
*Not appropriate for many research and production environments*

# Introducing Synthetic Experts

## What we don't always need

- Vast body of knowledge
- Ability to carry out many different tasks

*Sculpture by Jaume Plensa*

## What we often need

- Identify a **specific** (complex) construct
- No third-party constraints
- Accurate, efficient, reproducible

*just a tiny piece*

### Synthetic Experts

**Approximate** powerful Artificial Intelligence
*with an open-source Large Language Model (LLM)*

- Fine-tune pretrained LLM for classification task of interest
- Use powerful generative AI to label training data

*Pulling the Lollypop **out of the Candy Store***

# Empirical Example: The Marketing Mix on Social Media

**Construct of Interest: The Marketing Mix (MMX)**

- Marketing mix is **at the heart of marketing strategy**
- Theoretically founded
- 4 Ps: Product, Place, Price Promotion

**Data: Consumers' posts on Twitter mentioning brands**

- Vast and unstructured information source to marketers
- **Lens on hearts and minds of consumers**
- 3 years of Tweets for 699 brands

**Task: Identify MMX in 1,000 Tweets**

- **Human Experts**          *four academics*

  **vs.**

- **Crowdsourced Amateurs**  *from Amazon mTurk*
- **GPT-4**                  *via OpenAI API*
- **Synthetic Expert**       *trained on 15K Tweets labeled by GPT-4 \**

*@Sony's XM3's ain't as sweet as my bro's airpod pros but got a real steal 🤑 the other day #deal #headphonez*

*@dominos the other nite. waited over 1hr 4 cold pizza! SMH what's up with that?*

*I wish @abercrombie would stop using #usps to deliver their goods on this occasion, they give an email and text stating delivery between a 4 hour period. This is the 3rd delivery recently where I've been in all day waiting and nothing has arrived 😡*
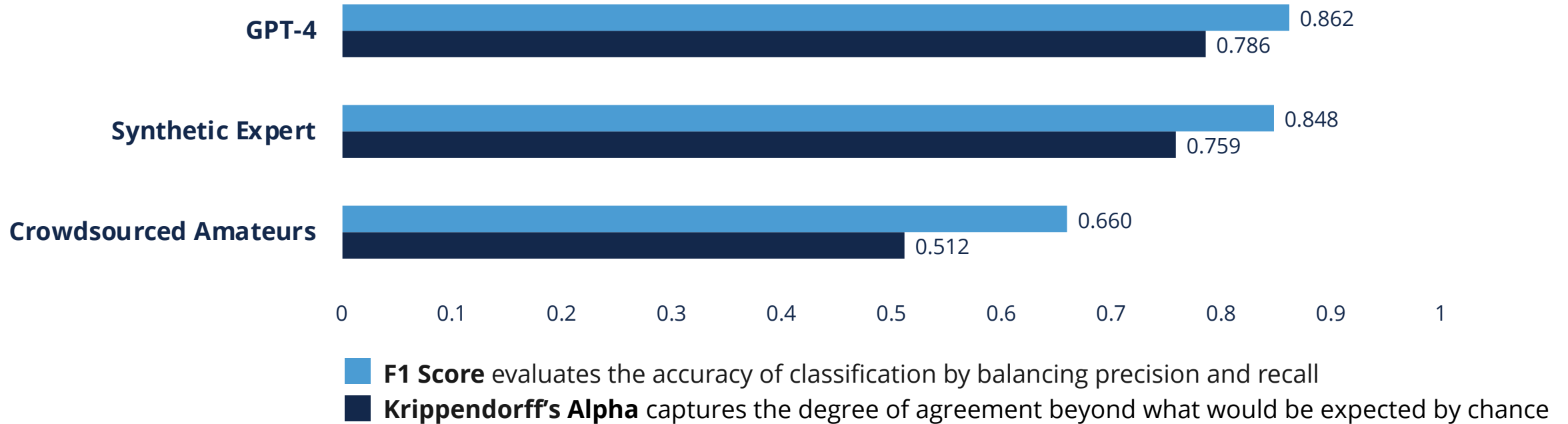
*Best cushioning ever!!! 🤗🤗🤗 my zoom vomeros are the bomb🏃💨!!! @nike #run #training*

*\* one-time cost of $90 to label Tweets with GPT-4, 30min fine-tuning of pre-trained open-source LLM RoBERTa on MacBook Pro*

# Main Findings

## Agreement with Human Expert Labels



| | F1 Score | Krippendorff's Alpha |
|---|---|---|
| GPT-4 | 0.862 | 0.786 |
| Synthetic Expert | 0.848 | 0.759 |
| Crowdsourced Amateurs | 0.660 | 0.512 |

■ **F1 Score** evaluates the accuracy of classification by balancing precision and recall
■ **Krippendorff's Alpha** captures the degree of agreement beyond what would be expected by chance

1. ***Crowdsourced labels insufficient*** for complex constructs
2. ***GPT-4 is viable surrogate*** for human expertise
3. ***Synthetic Expert*** is only ***2.5% less accurate,*** but ***66 x faster*** and ***400 x cheaper*** *than GPT-4\**

*\* on a standard MacBook Pro*

# Multiple Classifiers for richer Insight
## Brand Sentiment x MMX

| | Brand | Sentiment Overall | Sentiment by Marketing Mix Variable | | | |
|---|---|---|---|---|---|---|
| | | | Product | Place | Price | Promotion |
| **Apparel** | Calvin Klein | .256 | .258 | .339 | .419 | .330 |
| | Abercrombie & Fitch | .286 | .302 | .209 | .347 | .494 |
| | Polo Ralph Lauren | .291 | .295 | .336 | .192 | .408 |
| **Snacks** | Tostitos | .245 | .243 | .303 | .135 | .410 |
| | Doritos | .261 | .258 | .148 | .317 | .438 |
| | SunChips | .188 | .189 | .295 | .486 | .275 |
| **Banks** | Spirit Airlines | .024 | .002 | .025 | .072 | .532 |
| | JetBlue | .237 | .264 | .193 | .119 | .541 |
| | Southwest Airlines | .280 | .301 | .279 | .200 | .571 |

**Legend** — lowest | **Color indicates Relative Sentiment *** | highest

* column-wise for overall sentiment; row-wise across marketing mix variables

***Notes***: VADER Sentiment Analysis (Hutto and Gilbert 2014), 9 Brands, N = 9,000 randomly sampled Tweets from 2020; stratified by brand

# Implications

## Organizations

**Powerful, scalable, and accessible** solution for complex classification tasks
- free of third-party constraints
- mitigate privacy and confidentiality concerns
- easily updated or replaced when world and/or task changes

## Research

**Inform downstream tasks as DV or IV**
- access to complex constructs
- answer research questions
- test hypotheses
- explain mechanisms

**Ensure replicability** of research that leverages AI

*Hugging Face Model Hub*
*> 2,100 sentiment classifiers*

*NEW: 1 MMX Classifier*
***66 x faster*** *and*
***400 x cheaper*** *than GPT-4*

# Outlook

## Many sectors and functions

- **Marketers** may investigate *service quality* dimensions, *customer experience* dimensions, *branding* (identity, equity, image), or elements for a *SWOT* analysis.

- **Lawyers** may seek *contract elements* such as offer, acceptance, and consideration in memoranda, addendums, and communications.

- **Policy makers** may need to identify *agenda-setting and policy frames* in documents, government communications, or news reports.

- **Organizations** may want to understand *leadership styles* such as autocratic, democratic, or laissez-faire from corporate communications, internal memos, or employee reviews.

## Future innovations in AI ...

- **improve performance** of Synthetic Experts further

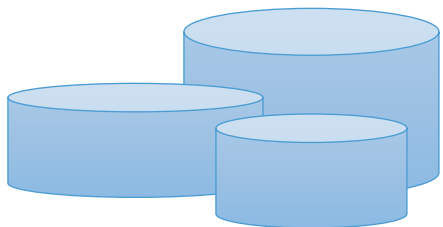- promise to **extend** Synthetic Experts across **different media** types

# Data Sharing

## Anonymizing texts with **Synthetic Twins**

**Basic Idea:** Use generative AI to Create replicas of texts that reflect their idea and meaning but obfuscate identifying information

### Synthetic Twins

- Correspond semantically in idea and meaning to original texts

- Wording, people, places, firms, brands, and products changed by AI

- Mitigate, to some extent, possible privacy, and copyright concerns

- Can be useful to create variations of existing texts

**Multiple Demo Datasets available _here_**
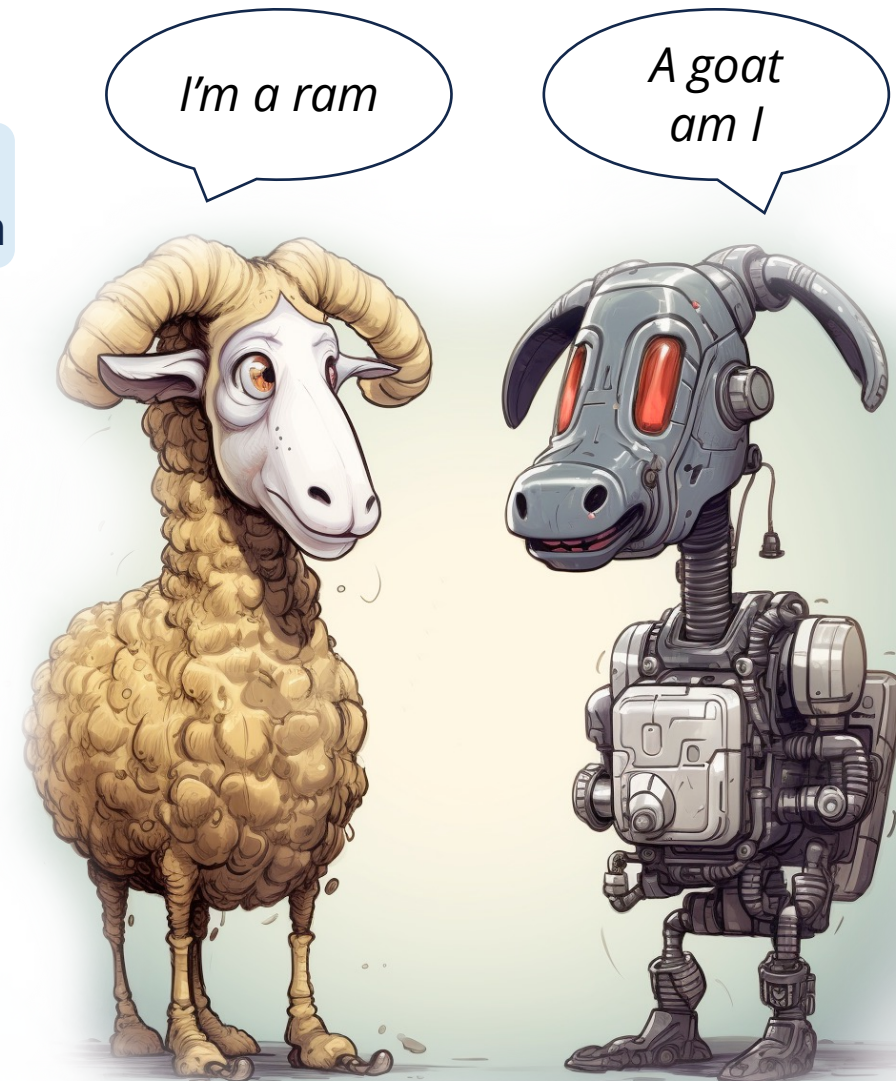
**Create your own Synthetic Twins _here_**

*I'm a ram*

*A goat am I*

*Image by Midjourney*

# Knowledge Sharing: www.synthetic-experts.ai

## *Website*     *Code*     *Working Paper*

### How to Create your own Synthetic Expert using Generative AI

**This website provides supporting materials, code, data, and tutorials for the paper** *Creating Synthetic Experts with Generative AI* **by Daniel M. Ringel (2023).**

#### News

**[October 19, 2023] Recording of Wharton AI Talk**: My talk on Synthetic Experts is now available on YouTube

**[September 28, 2023] Top 10 at SSRN**: Synthetic Experts **rank #1** in Top 10 downloads of recent papers in Marketing Science eJournal Ranking

**[September 8, 2023]** Presentation of *Synthetic Experts* at the Wharton Business & Generative AI Workshop in San Fransico. Check out the PDF presentation slides!

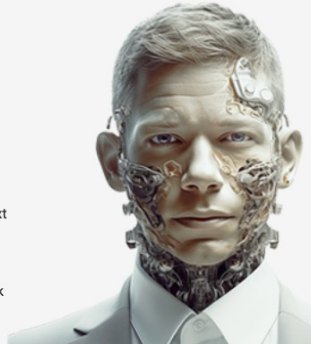**[September 6, 2023]** Anonymize texts with *Synthetic Twins*: Python notebook now available on GitHub

**[August 26, 2023]** Create your own *Synthetic Expert* with new Python notebooks available on GitHub. Label text with generative AI and fine-tune LLMs for your purpose.

**[August 23, 2023]** Work with several demo Datasets to explore the use and creation of **Synthetic Experts**.

**[August 19, 2023]** Easily apply the **MMX Synthetic Expert** of this research to YOUR data. New Python notebook available on GitHub.

**[August 16, 2023]** *This research is currently being revised and extendet. I will make materials available as they are completed.*

#### Paper

- Ringel, Daniel, *Creating Synthetic Experts with Generative AI* (July 15, 2023). Available at SSRN: https://ssrn.com/abstract=4542949.
- Appendix: Details, Notes, Parameters, IRB information

#### Code

- Python Notebooks: GitHub Repo
- Functions: UseSynExp.py

#### Data

- List of Brands
- List of Tweets
- Synthetic Twins of Data

#### Supporting Documents

- How to set-up GPU computing on Apple Silicone [ipynb]
- How to set-up your own Deep Learning Machine: Install Ubuntu with CUDA, CuDNN and PyTorch on a computer [pdf]
- How to run the code on Google Colaboratory (for free) [pdf]
- **BONUS**: Synthetic Twins of real-world textual data [ipynb]

#### Tutorials

- Get ready to fine-tune LLMs [Setup_Python_GPU.ipynb]
- Query the OpenAI API [Query_OpenAI_API.ipynb]
- Parse OpenAI API responses [Parsing_API_Responses.ipynb]
- Fine-tune a pre-trained LLM from Huggingface [Fine-tuning_LLMs.ipynb]
- Hyperparametertuning with Optuna [Hyperparametertuning.ipynb]
- Deploy your Synthetic Expert and use it for Inference [ipynb]

#### Classroom Materials

- Slides [pptx]
- Notebook [ipynb]
- Data [zip]
- Instructor Notes [pdf]
- Video [YouTube]

#### Repository

Code will be maintained on GitHub

#### Fine-Tuned Model

Marketing Mix Classifier on Huggingface's Model Hub

#### Questions, Comments?

Get in touch at dmr@unc.edu

## *MMX Synthetic Expert*



## *contact* dmr@unc.edu