# Who Am I?

**Jacob Marks**

B.S. in Math & Physics @ Yale

Ph.D. Theoretical Physics @ Stanford

ML Engineer & Developer Evangelist at Voxel51

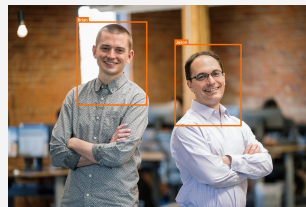Follow me on LinkedIn!

voxel51.com

# Who We Are

## Voxel51

The lead maintainers of the open source FiftyOne toolset

## Innovators

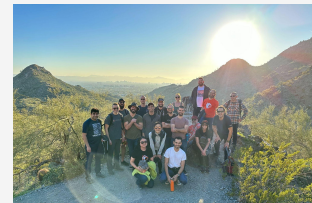Spun out of the Computer Vision Lab at the University of Michigan

## Our Founders

**Jason Corso, PhD**
Professor of Computer Vision at the University of Michigan + 20 years working in CV and ML

**Brian Moore, PhD**
Computer Vision Extraordinaire + Built unique algorithms for CV at University of Michigan + Deep Expert in ML

## Focused

**Team of 25+** driven, sharp, talented Computer Scientists and CV Experts supporting thousands of open source FiftyOne users and Fortune 100 companies

voxel51.com

# Voxel51's Mission

is to bring transparency and clarity to the world's data

## Data Quality

Nothing hinders machine learning systems more than **poor quality data**

The secret to getting ML products into production lies in **improving datasets, not model architectures**

## Data-Centric Workflows

Companies are hiring ML talent as fast as they can, but they **need tools and workflows** to build high quality datasets and models

Our **open source approach** aims to bring data-centric ML to anyone who wants it

# What Is FiftyOne?

**In a Nutshell**

> Visualize, clean, and curate

> Find hidden structure

> Evaluate models

> *Flexible. Customizable. Connected.*

***An open source tool for building high-quality datasets and computer vision models***

voxel51.com

# Agenda

> **LLMs and RAG**

> Going Multimodal

> Testing Multimodal RAG Pipelines

> Next Steps

voxel51.com

# LLMs Explode

# LLMs Explode

# How LLMs Work



🔗 https://bbycroft.net/llm

voxel51.com

9

# How LLMs Work



The sky is → LLM ┈┈┈

blue = -0.96
clear = -1.60
usually = -2.47
the = -3.40
< = -3.47

More likely
Less likely

→ The sky is blue

Total: -0.96 logprob on 1 token
(73.18% probability covered in top 5 logits)

# LLM Limitations



> Knowledge cutoff

> Domain-specific tasks

> Hallucination

# Retrieval Augmented Generation (RAG)

# RAG Modifies LLM Inputs

# How We Use RAG at Voxel51

# Decomposing RAG

**Data Sources**

**Vector DB**

**Technique**

**LLM Framework**

voxel51.com

# Agenda

> LLMs and RAG

> **Going Multimodal**

> Testing Multimodal RAG Pipelines

> Next Steps

voxel51.com

# GPT-4 (All Tools)



voxel51.com

17

# Adept's Fuyu-8b



The HBO Recycling Program
Actors who have appeared in three or more episodes of multiple scripted, live-action, original HBO series since Oz (excluding miniseries)

* Apologies that J.D. Williams is out of alphabetical order.
Research: Sam Schube and Danny Savitzky
Visualization: Craig Robinson

**Question:** "Aidan Gillen acted in how many series?"
**Fuyu's answer:** "2"



**Question:** "What is the highest life expectancy at birth of males?"

**Fuyu's answer:** "The life expectancy at birth of males in 2018 is 80.7"

# Multimodal LLM Landscape



GPT-4 Vision



CogVLM 更强的 VisualGLM



Qwen-VL



Gemini



MiniGPT-4



Adept Fuyu 8b

# Multimodal LLMs Accept *Multiple* Images

# Applications of Multimodal LLMs: Medical

Med-PaLM

# Applications of Multimodal LLMs: Retail

# Multimodal LLMs Have Similar Limitations



https://arxiv.org/pdf/2302.04023.pdf

# Multimodal Data Helps Humans Learn

voxel51.com

# Multimodal RAG

# Agenda

> LLMs and RAG

> Going Multimodal

> **Testing Multimodal RAG Pipelines**

> Next Steps

voxel51.com

# Quantitative Evaluation for Multimodal RAG

**Retrieval Eval**

> *Hit Rate*

> *MRR*

⭐ ~Separately for image and text~

**Generation Eval**

> *Correctness*

> *Faithfulness*

> *Relevancy*

🔗 https://docs.llamaindex.ai/en/stable/examples/evaluation/multi_modal/multi_modal_rag_evaluation.html

# FiftyOne's Multimodal RAG Plugin!



```
> pip install fiftyone
>
> fiftyone plugins download https://github.com/jacobmarks/fiftyone-multimodal-rag-plugin
> fiftyone plugins requirements @jacobmarks/multimodal_rag --install
```

LIVE DEMO

# Agenda

> Data, Data, Data

> The ABCs of Multimodal Search

> 4 Performance Enhancers

> **Next Steps**

voxel51.com

# Thank You!

# Join the Community



**Star us**



**Follow us**



**Try FiftyOne**

voxel51.com