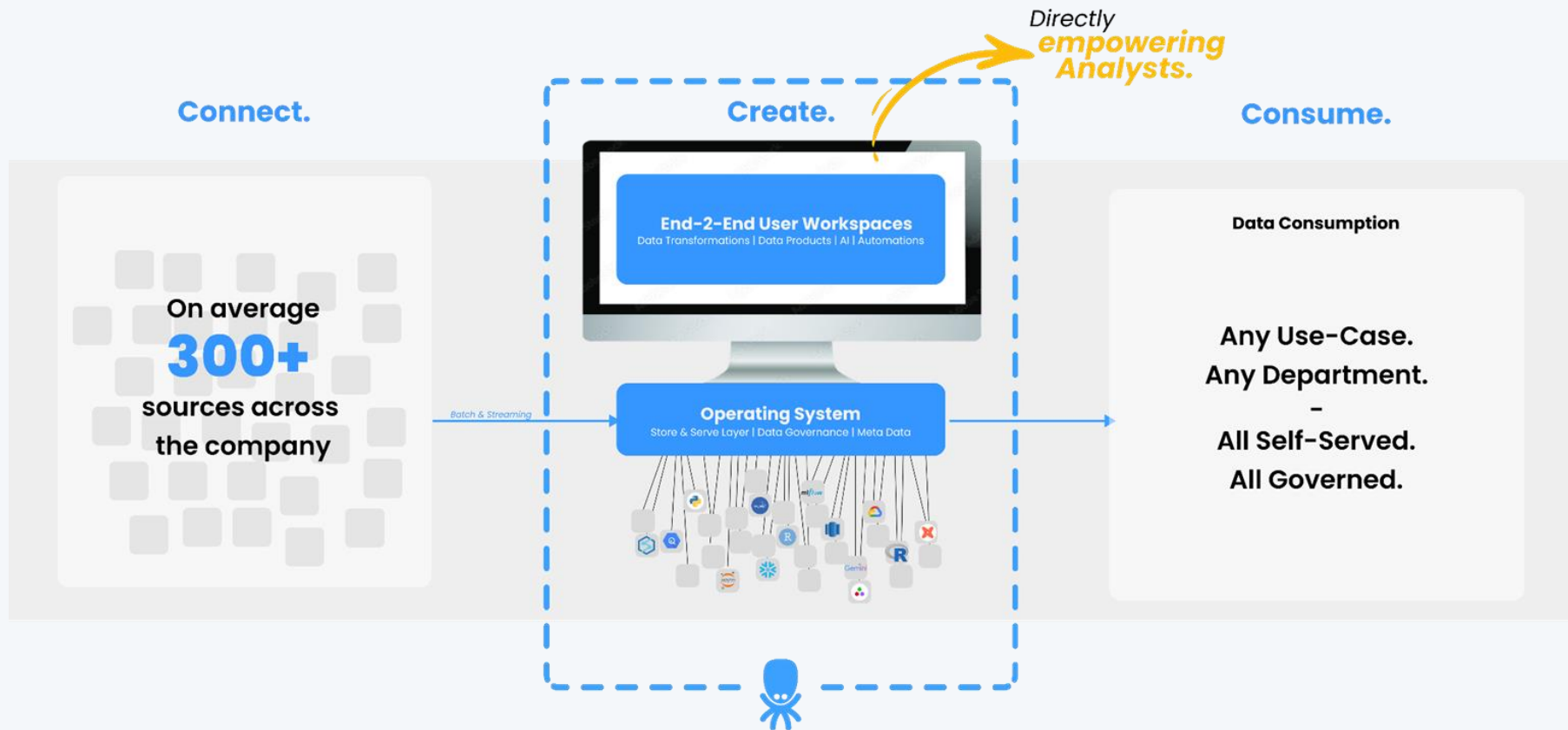# From Truisms to Transformation:

## Achieving Tangible Results in AI and Data Management

**Jordan Burger**
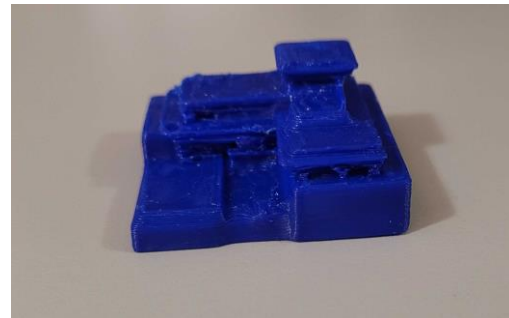Head of AI Lab, Keboola

# Unlock the **full potential of enterprise data.**

Keboola

*Directly* **empowering Analysts.**

**Connect.**

**Create.**

**Consume.**

On average **300+** sources across the company

**End-2-End User Workspaces**
Data Transformations | Data Products | AI | Automations

**Data Consumption**

**Operating System**
Store & Serve Layer | Data Governance | Meta Data

*Batch & Streaming*

**Any Use-Case.**
**Any Department.**
**-**
**All Self-Served.**
**All Governed.**

# **About** Me

- Research Lead; AI Applications

- Exploring and implementing AI technologies in data pipelines

- Addressing GenAI disruption in data and analytics



'A house in the style of Frank Lloyd Wright's Fallingwater'

Generated with Point-E from OpenAI

# ... But it's still right twice a day

Keboola

| Garbage **in** | > | Gospel **Out** |

# Struggling with Outliers;
Unable to handle genuine novelty

How many times does the letter "i" appear in "artificial"?

The letter "i" appears twice in the word "artificial."

**Keboola**

## Struggling with Outliers;

Unable to handle genuine novelty

**You**
A man and his mother are in a car accident. The mother sadly dies. The man is rushed to the ER. When the doctor sees him, he says, "I can't operate on this man. He's my son!"
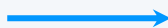
How is this possible?

**ChatGPT**
The doctor is the man's other parent—his mother, indicating that the doctor is a woman. This riddle plays on common assumptions about professions and gender roles.
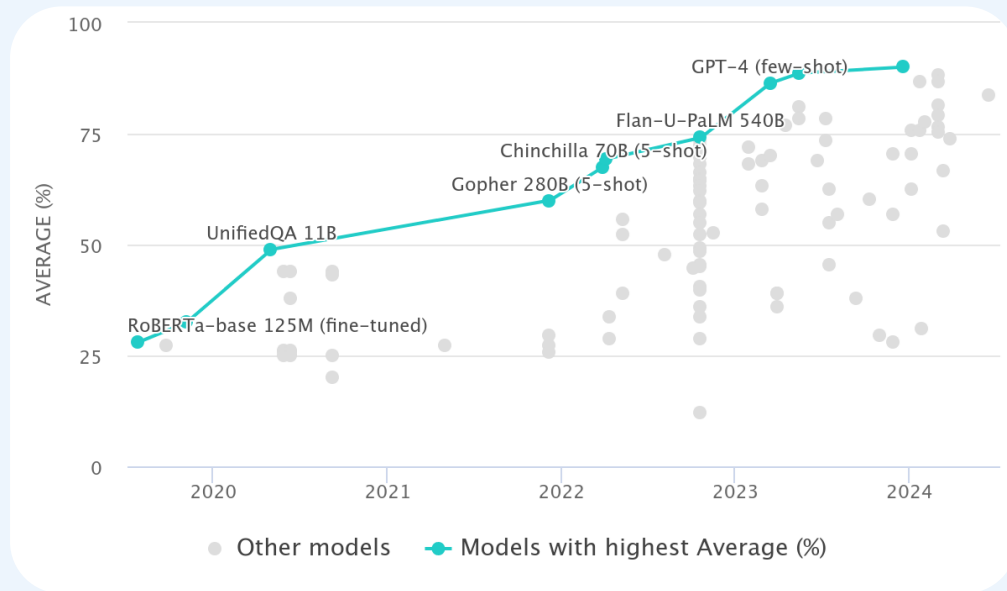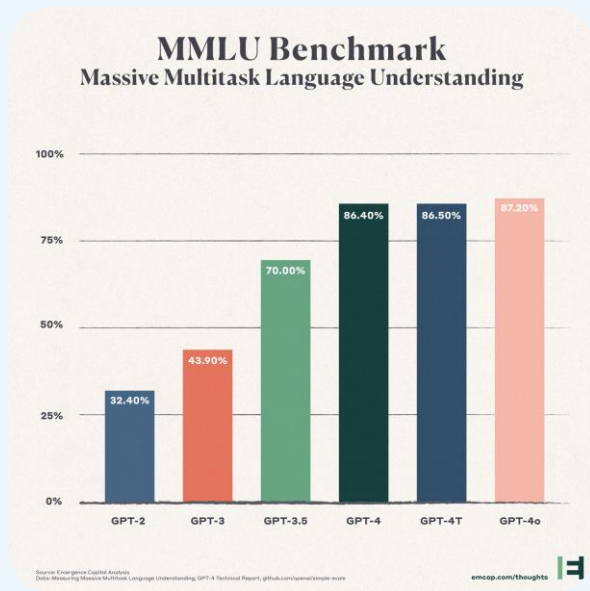
Credit: @colin_fraser

# Struggling with Outliers;
## Unable to handle genuine novelty

| Employee Salaries | |
|---|---|
| Employee ID | Salary |
| 12 | 500 CZK/Month |
| 116 | 1.000.000$/year |
| 1243 | $75,000/year |
| 242344 | 3.000/month |
| 555 | 250.00 Kc/Year |
| 63465 | 100.000 USD/month |

| Employee Salaries | | | |
|---|---|---|---|
| Employee ID | Currency | Salary | Pay Frequency |
| 12 | USD | $5,000 | Annual |
| 116 | USD | $1,000,000 | Annual |
| 1243 | USD | $75,000 | Annual |
| 242344 | USD | $360,000 | Annual |
| 555 | USD | $32,000 | Annual |
| 63465 | USD | 1,200,00 | Annual |

# **The Plateau** was always (kind of) here



MMLU Benchmark
Massive Multitask Language Understanding

GPT-2: 32.40%
GPT-3: 43.90%
GPT-3.5: 70.00%
GPT-4: 86.40%
GPT-4T: 86.50%
GPT-4o: 87.20%

Source: Emergence Capital Analysis
Data: Measuring Massive Multitask Language Understanding, GPT-4 Technical Report, github.com/openai/simple-evals
emcap.com/thoughts



AVERAGE (%)

GPT-4 (few-shot)
Flan-U-PaLM 540B
Chinchilla 70B (5-shot)
Gopher 280B (5-shot)
UnifiedQA 11B
RoBERTa-base 125M (fine-tuned)

● Other models   ●— Models with highest Average (%)

# The Machines Are Taking Over

**COMPUTERS OUTDO MAN AT HIS WORK NOW—AND SOON MAY OUTTHINK HIM**

COMPUTER'S INNARDS are examined by one of humans who tend it. IBM programmer looks through glass panel in door that gives access to heat control and memory centers.
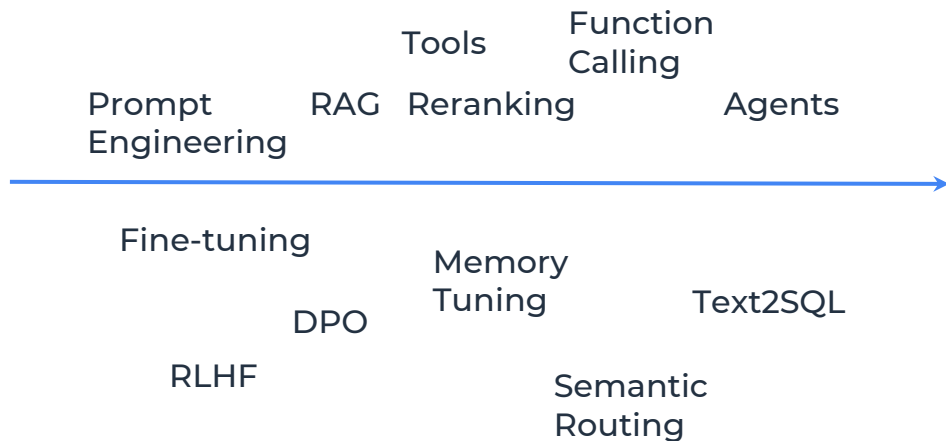
# But where's the value?

"

**Investment in AI has reached a new high with a focus on generative AI, which, in most cases, has yet to deliver its anticipated business value.**

**Gartner – Hype Cycle for Artificial Intelligence, 2024**

Prompt Engineering    RAG    Reranking    Tools    Function Calling    Agents

Fine-tuning

DPO    Memory Tuning    Text2SQL

RLHF    Semantic Routing

# But where's the value?

> The revenue isn't there yet, and might never come.

**Gary Marcus**

All of the applied research going on can be boiled down to one thing: "These models are cool; now how can I get them to do what I actually need?"

# *Molmo and PixMo*:
## Open Weights and Open Data
## for State-of-the-Art Multimodal Models

Matt Deitke*[†][ψ]    Christopher Clark*[†]    Sangho Lee[†]    Rohun Tripathi[†]    Yue Yang
Jae Sung Park[ψ]    Mohammadreza Salehi[ψ]    Niklas Muennighoff[†]    Kyle Lo[†]    Luca Soldaini[†]
Jiasen Lu[†]    Taira Anderson[†]    Erin Bransom[†]    Kiana Ehsani[†]    Huong Ngo[†]
YenSung Chen[†]    Ajay Patel[†]    Mark Yatskar[†]    Chris Callison-Burch[†]    Andrew Head
Rose Hendrix[†]    Favyen Bastani[†]    Eli VanderBilt[†]    Nathan Lambert[†]    Yvonne Chou[†]
Arnavi Chheda[†]    Jenna Sparks[†]    Sam Skjonsberg[†]    Michael Schmitz[†]    Aaron Sarnat[†]
Byron Bischoff[†]    Pete Walsh[†]    Chris Newell[†]    Piper Wolters[†]    Tanmay Gupta[†]    Kuo-Hao Zeng[†]
Jon Borchardt[†]    Dirk Groeneveld[†]    Jen Dumas[†]    Crystal Nam[†]    Sophie Lebrecht[†]
Caitlin Wittlif[†]    Carissa Schoenick[†]    Oscar Michel[†]    Ranjay Krishna[†][ψ]    Luca Weihs[†]
Noah A. Smith[†][ψ]    Hannaneh Hajishirzi[†][ψ]    Ross Girshick[†][ψ]    Ali Farhadi[†][ψ]    Aniruddha Kembhavi[†][ψ]

[†]Allen Institute for AI    [ψ]University of Washington

prietary. **The strongest open-weight models rely heavily on synthetic data from proprietary VLMs to achieve good performance, effectively distilling these closed models into open ones.** As a result, the community is still missing foun-

are state-of-the-art in their class of openness. *Our key innovation is a novel, highly detailed image caption dataset col-*

data. *The success of our approach relies* on careful choices for the model architecture details, a well-tuned training pipeline, and, *most critically, the quality of our newly collected datasets,* all of which will be released. The best-in-

## 1. Introduction

Extensions to large language models (LLMs) that process images in addition to text have resulted in impressive multimodal capabilities, such as generating comprehensive image descriptions and accurately answering complex visual questions. The most performant of these vision-language models, however, remain proprietary with neither model weights, data, nor code being publicly released.

With the goal of fostering scientific exploration, numerous research efforts have attempted to reproduce similar capabilities in *open* models. Early works, exemplified by LLaVA [15], produced fully open weights and training data but now lag significantly behind the state-of-the-art. More recent, stronger open-weight models have trended towards less open data: the training data may either be proprietary (*e.g.*, [5]) or, in cases where it is released, there is a heavy reliance on *synthetic* data generated by proprietary systems, *e.g.*, models are trained on datasets like ShareGPT4V [7] which uses GPT-4V [25] to generate a large set of detailed

nect an independently pre-trained, off-the-shelf vision encoder and language model and jointly train the resulting VLM to generate captions from a newly collected dataset of detailed, high-quality, dense image descriptions. After joint training, we follow standard practice and use supervised fine-tuning to produce an instruction following model.

quality data (*e.g.*, [4, 5]). The success of our approach relies on careful choices for the model architecture details, a well-tuned training pipeline, and most critically, the quality of our new datasets, collectively named **PixMo** (**Pix**els for **Mo**lmo), all of which will be released.

well-tuned training pipeline, and most critically, the quality of our new datasets, collectively named **PixMo** (**Pix**els for **Mo**lmo), all of which will be released.

In practice, it is challenging to collect dense captioning datasets from human annotators. If asked to write an image description, the result often only mentions a few salient visual elements [8]. If a minimum word count is enforced, annotators will either take too long to type, making collec-

data from proprietary VLMs. Our key innovation is a simple but effective data collection strategy that avoids these problems: we ask annotators to describe images in *speech* for 60 to 90 seconds rather than asking them to write descriptions. We prompt the annotators to describe everything

problems: we ask annotators to describe images in *speech* for 60 to 90 seconds rather than asking them to write descriptions. We prompt the annotators to describe everything they see in great detail, including descriptions of spatial positioning and relationships. Empirically, we found that with this modality switching "trick" annotators provide far more detailed descriptions in less time, and for each description, we collect an audio receipt (*i.e.*, the annotator's recording) proving that a VLM was not used.

After training our models to generate dense captions we



Figure 1. The **Molmo** architecture follows the simple and standard design of combining a language model with a vision encoder. Its strong performance is the result of a well-tuned training pipeline and our new **PixMo** data.

us to rank models by user preference. Our smallest model, MolmoE-1B, based on the OLMoE-1B-7B mixture-of-experts LLM, nearly matches the performance of GPT-4V on both academic benchmarks and user preference. Molmo-7B-O and Molmo-7B-D, based on OLMo-7B [10] and Qwen2 7B [33], respectively, perform comfortably between GPT-4V and GPT-4o on both academic benchmarks and user preference. Our best-in-class Molmo-72B model, based on Qwen2 72B, achieves the highest academic benchmark score and ranks second by hu-

# NVLM: Open Frontier-Class Multimodal LLMs

Wenliang Dai*     Nayeon Lee*     Boxin Wang*     Zhuoling Yang*

Zihan Liu   Jon Barker   Tuomas Rintamaki   Mohammad Shoeybi   Bryan Catanzaro

Wei Ping*,†

**NVIDIA**

* {wdai, nayeonl, boxinw, zhuoliny, wping}@nvidia.com

## Abstract

We introduce NVLM 1.0, [1] a family of frontier-class multimodal large language models (LLMs) that achieve state-of-the-art results on vision-language tasks, rivaling the leading proprietary models (e.g., GPT-4o) and open-access models (e.g., Llama 3-V 405B and InternVL 2). Remarkably, NVLM 1.0 shows improved text-only performance over its LLM backbone after multimodal training.

In terms of **model design**, we perform a comprehensive comparison between decoder-only multimodal LLMs (e.g., LLaVA) and cross-attention-based models (e.g., Flamingo). Based on the strengths and weaknesses of both approaches, we propose a novel architecture that enhances both training efficiency and multimodal

and supervised fine-tuning datasets. Our findings indicate that dataset quality and task diversity are more important than scale, even during the pretraining phase, across all architectures. Notably, we develop **production-grade multimodality**

## 3.1 Data

The AFM pre-training dataset consists of a diverse and high quality data mixture. This includes data we have licensed from publishers, curated publicly-available or open-sourced datasets, and publicly available information crawled by our web-crawler, Applebot [Apple, 2024a]. We respect the right of webpages to opt out of being crawled by Applebot, using standard robots.txt directives

Given our focus on protecting user privacy, we note that no private Apple user data is included in the data mixture. Additionally, extensive efforts have been made to exclude profanity, unsafe material, and personally identifiable information from publicly available data (see Section 7 for more details). Rigorous decontamination is also performed against many common evaluation benchmarks.

We find that data quality, much more so than quantity, is the key determining factor of downstream model performance. In the following, we provide more details about key components of the data mixture.

### 3.1.1 Web pages

We crawl publicly available information using our web crawler, Applebot [Apple, 2024a], and respect the rights of web publishers to opt out of Applebot using standard robots.txt directives. Plus, we take steps to exclude pages containing profanity and apply filters to remove certain categories of personally identifiable information (PII). The remaining documents are then processed by a pipeline

**Keboola**

# Data is Central to AI Training
Foundation Models

## We will run out of data? An analysis of the limits of scaling datasets in Machine Learning

Projections of Low-Quality Data

Keboola

**Data is Central to AI Training** Foundation Models

Will we run out of data? An analysis of the limits of scaling datasets in Machine Learning

Projections of High-Quality Data

**Keboola**

**What are you doing to lengthen this curve in *your* domain?**

Will we run out of data? An analysis of the limits of scaling datasets in Machine Learning

Projections of High-Quality Data

## Job Error

**SQL Error: Type mismatch: Cannot concatenate a string with an array**
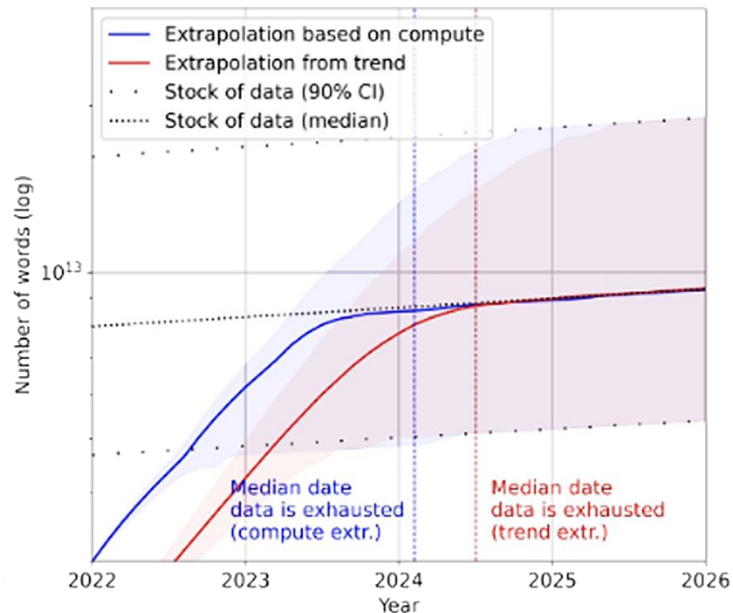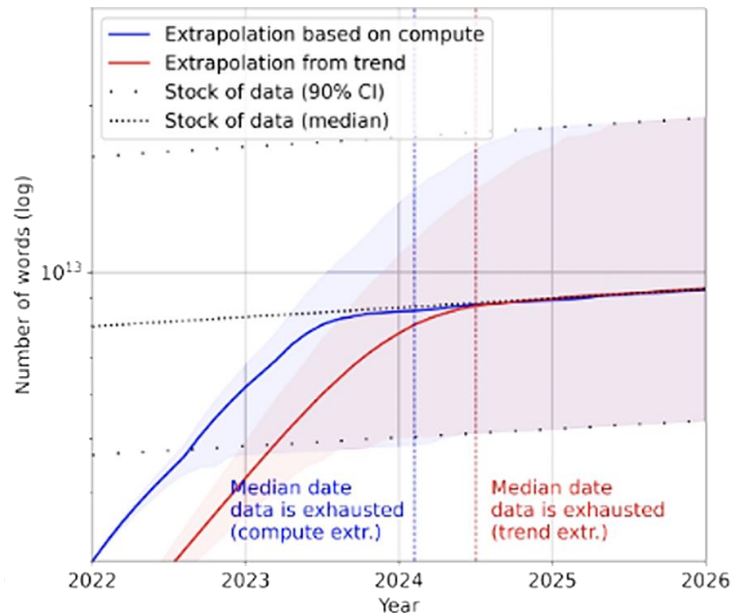
### AI Explanation

This error is caused by trying to concatenate a string with an array. **You should first convert the array to a string and try again.** You can use the ARRAY_TO_STRING function to resolve this issue.

Is this AI-generated explanation useful?

COPY DETAILS TO CLIPBOARD

---

## Job Error

**SQL Error: Type mismatch: Cannot concatenate a string with an array**

### AI Explanation

This error occurs because the SQL query attempts to concatenate a string with an array, which is not supported. **The root cause is a deliberate change made recently in an upstream dependency, which altered the data type of the column being concatenated.**

**Root Cause Analysis:**

- **Lineage Tracking:** This error traces back to a recent update in the user_data_transformation step, where the user_interests column was changed from a string to an array. This change affects the current query in the report_generation_transformation step, where the user_interests column is concatenated with other strings.

Is this AI-generated explanation useful?

COPY DETAILS TO CLIPBOARD

CONTACT SUPPORT

## revenuecategory Table Description

This table contains data for tracking various revenue categories. Each category represents different segments related to fitness, retail, training gear, climbing accessories, and retail packs.

- Primary key: The table is uniquely identified by the id column.
- Columns:
  - id: A unique identifier for each revenue category.
  - name: The name of the revenue category.
  - description: A detailed description of the revenue category (e.g. "Fitness, Retail - Training Gear, Retail - Climbing Accessories, etc.").
  - isActive: A boolean
  - createdAt: The time
  - updatedAt: The tim
  - createdBy: Metada
  - updatedBy: Metada
  - giftCardTypeId: Ad

🧠 Generate Description

Is this AI-generated desc

---

## revenuecategory Table Description

CANCEL · ✓ SAVE

The revenuecategory table is central to the forecasting and business insights process within the Keboola project. It contains categorized revenue data across product lines such as Fitness, Retail, and Climbing Accessories. This classification differentiates revenue streams by location and type, enabling detailed forecasting models.

In the context of the Global Historical Revenue Forecasting transformation, this table is used for the multivariate revenue forecast, where each revenue entry is tied to a specific revenue category and location. By using the data from revenuecategory, the project generates segmented revenue predictions, helping the business understand how different categories perform across regions.

This segmentation enables the business to derive insights, such as identifying which categories drive revenue growth in specific locations or identifying underperforming segments. These forecasts inform decisions on resource allocation, product focus, and regional strategies, based on how revenue categories are expected to perform in the future.
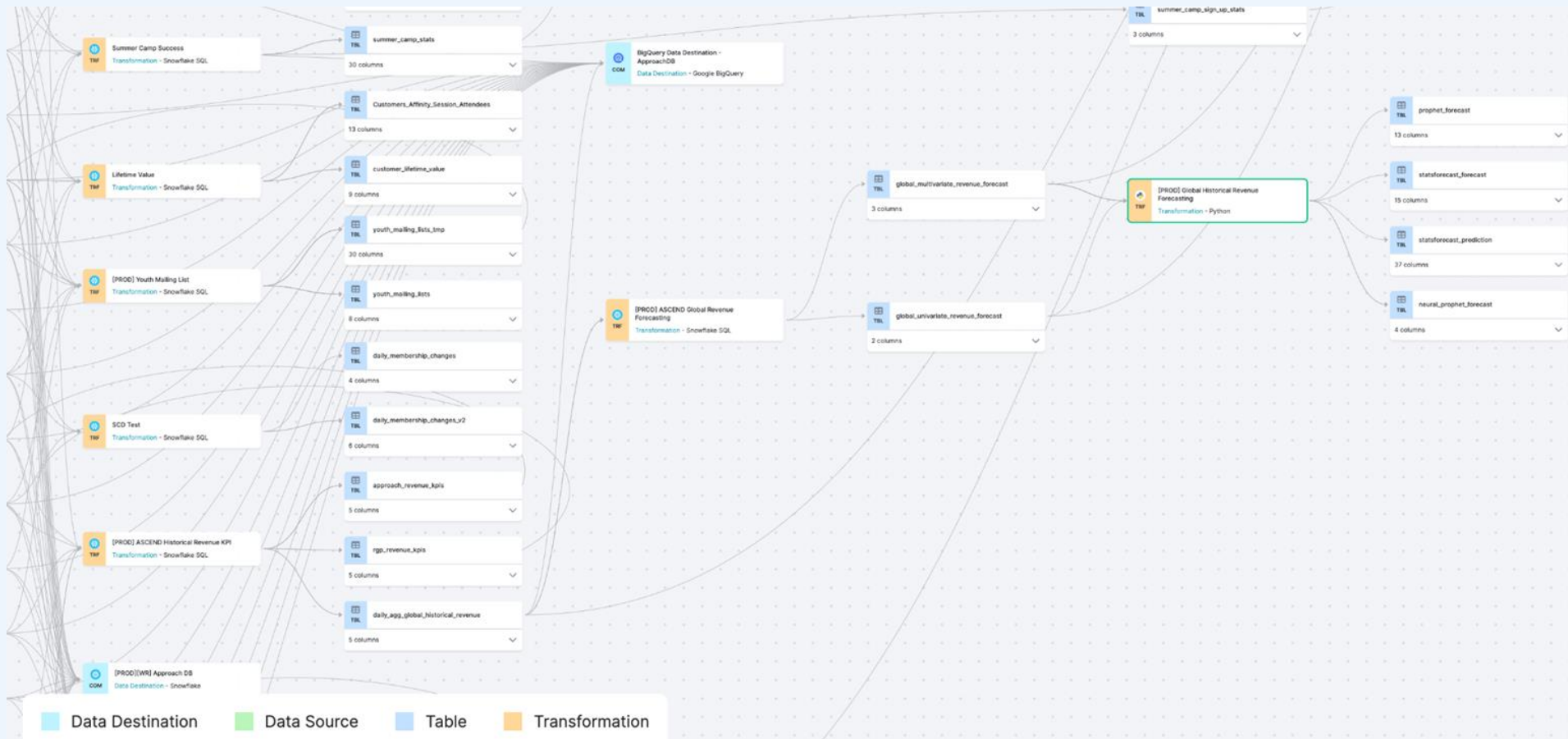
🧠 Generate Description                    Markdown is supported

Is AI-generated description useful?        👍 👎 | ↻

Summer Camp Success
Transformation - Snowflake SQL

summer_camp_stats
30 columns

Customers_Affinity_Session_Attendees
13 columns

BigQuery Data Destination - ApproachDB
Data Destination - Google BigQuery

summer_camp_sign_up_stats
3 columns

Lifetime Value
Transformation - Snowflake SQL

customer_lifetime_value
9 columns

prophet_forecast
13 columns

[PROD] Youth Mailing List
Transformation - Snowflake SQL

youth_mailing_lists_tmp
30 columns

global_multivariate_revenue_forecast
3 columns

[PROD] Global Historical Revenue Forecasting
Transformation - Python

statsforecast_forecast
15 columns

youth_mailing_lists
8 columns

[PROD] ASCEND Global Revenue Forecasting
Transformation - Snowflake SQL

statsforecast_prediction
37 columns

daily_membership_changes
4 columns

global_univariate_revenue_forecast
2 columns

neural_prophet_forecast
4 columns

SCD Test
Transformation - Snowflake SQL

daily_membership_changes_v2
6 columns

approach_revenue_kpis
5 columns

[PROD] ASCEND Historical Revenue KPI
Transformation - Snowflake SQL

rgp_revenue_kpis
5 columns

daily_agg_global_historical_revenue
5 columns

[PROD][WR] Approach DB
Data Destination - Snowflake

Data Destination   Data Source   Table   Transformation
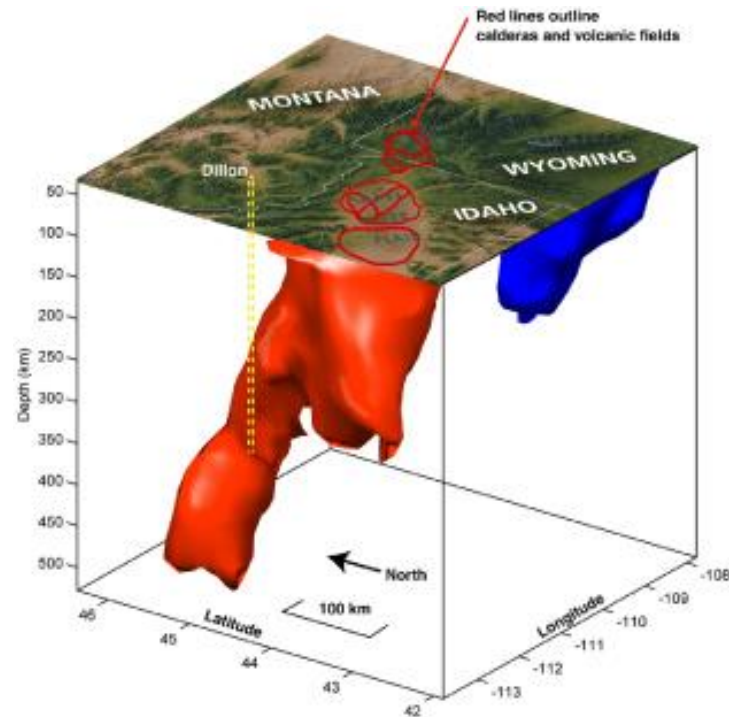
**Knowing the trajectory of AI doesn't help to predict the landscape that it will shape**

Keboola

# Building foundations for
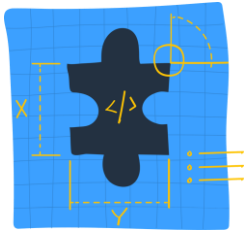AI is a delicate balance

# We need to demand **more** from our tools and vendors

- What is the data for? Does it represent what we want?

- What data is actually being used? How is it being processed?

# We need to demand **more** from our tools and vendors

- What is the data for? Does it represent what we want?

- What data is actually being used? How is it being processed?

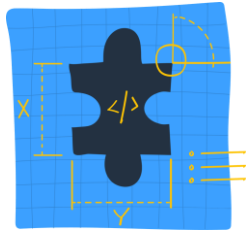- How do we validate that the data and models do the things that we want it to do?

Keboola

# Lessons from the Trenches



**Be Pragmatic**

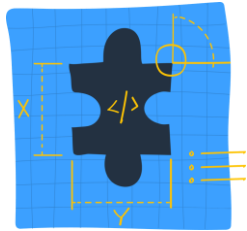# Lessons from the Trenches

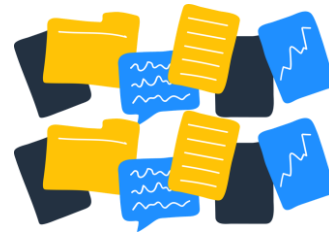### Be Pragmatic

### Forge Partnerships

# Lessons from the Trenches

**Be Pragmatic**

**Forge Partnerships**

**Curate your Knowledge**