# From Raw Data to Refined Datasets

Amir Alush, PhD

CTO @Visual Layer
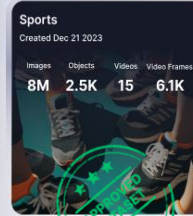
# Pixels to Products: My Tech Journey

- Ph.D. on Discrete Optimization Problems (in computer vision)

- Deep Learning since 2012

- Previously CTO & Co founder of **Brodmann17**

- Currently CTO & Co founder of **Visual Layer (fastdup)**

Visual Layer

# Family Album: 500GB of Chaos (and growing..)

- Married + 2 children + dog
- 500 hours of manually sifting



Generated by DALL·E 3

# Deja Vu at Work: 100TB of Chaos

- Automotive, tons of data collected on a daily basis

- All this data has to go through a curation process

**Raw Data**

**Dataset**

Visual Layer

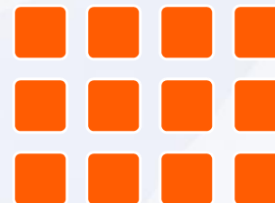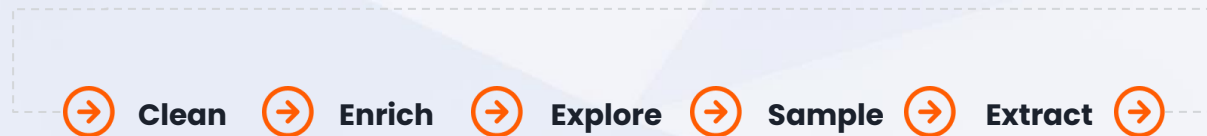# Deja Vu at Work: 100TB of Chaos

- Automotive, tons of data collected on a daily basis

- All this data has to go through a curation process

**Raw Data**

Dataset Curation Process

**Dataset**

→ **Clean** → **Enrich** → **Explore** → **Sample** → **Extract** →

Visual Layer

Dataset Curation Process today is mostly: **Done Manually** 👉

# The challenges of: Dataset Curation

**1**

**Speed**

Slow and Unscalable

**2**

**Cost**

Very expensive

**3**

**Quality**

Very poor quality
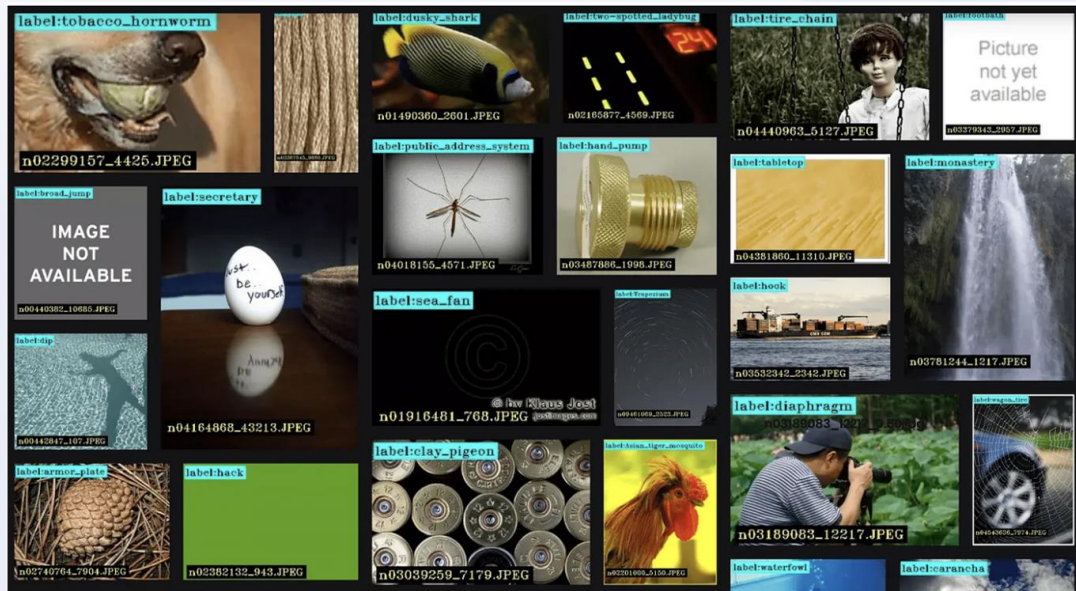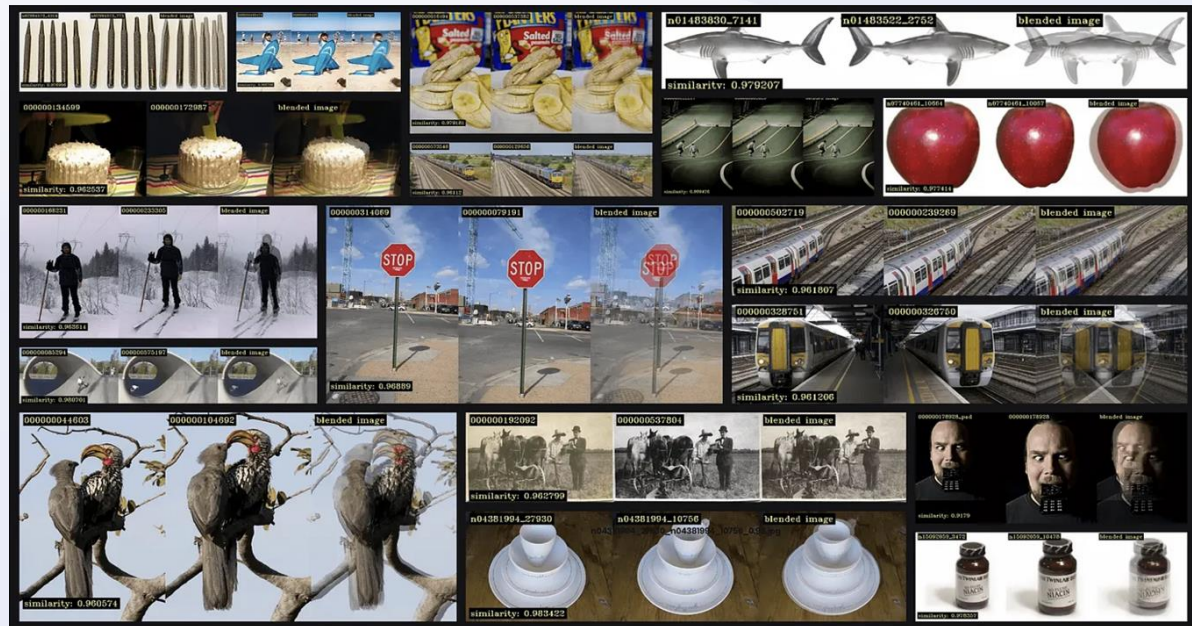
Visual Layer

# When Datasets Go Wrong: Bad Labels



Example wrong labels in the ImageNet-21K dataset. Additional information can be found in our GitHub repo.

Example of many different labels for the same object. There are thousands of such clusters.

large-image-datasets-today-are-a-mess

Visual Layer

# When Datasets Go Wrong: Duplicates



Example near duplicates identified in the MS-COCO (160K images) & ImageNet-21K datasets (11.5M images). A record breaking number of 1.2M duplicates were identified in the ImageNet-21K dataset! Additional

large-image-datasets-today-are-a-mess

Visual Layer

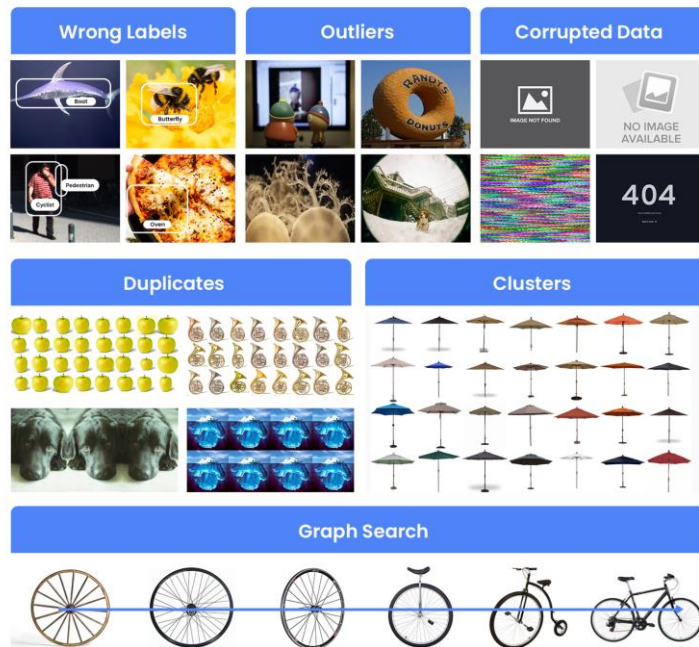# fastdup: Efficient Blocks for Dataset Curation

fastdup

**1.4K+**
GitHub stars

**380K+**
Downloads

**50B+**
Images processed

visual-layer/fastdup



Visual Layer

# fastdup: Efficient Blocks for Dataset Curation

In short, you'll need 3 lines of code to run fastdup:

```python
import fastdup
fd = fastdup.create(input_dir="IMAGE_FOLDER/")
fd.run()
```

## Enrich Data Using Foundation Models

The notebooks in this section show how to enrich your visual dataset using various foundation models supported in fastdup.

📚 **Zero-Shot Classification:** Enrich your visual data with zero-shot image classification and tagging models such as Recognize Anything Model, Tag2Text, and more.

🔗 Learn More.

🧭 **Zero-Shot Detection:** Enrich your visual data with zero-shot image detection model such as Grounding DINO and more.

🔗 Learn More.

## Load Data From Sources

The notebooks in this section show how to load data from various sources and analyze them with fastdup.

🤗 **Hugging Face Datasets:** Load and analyze datasets from Hugging Face Datasets. Perfect if you already have a dataset hosted on Hugging Face hub.

🔗 Learn More.

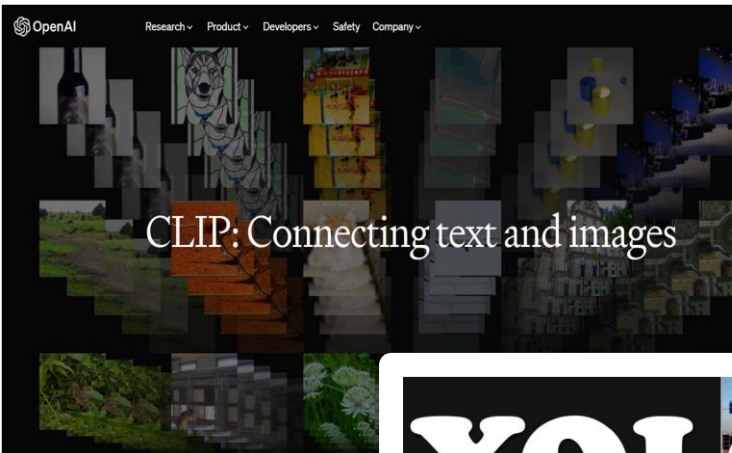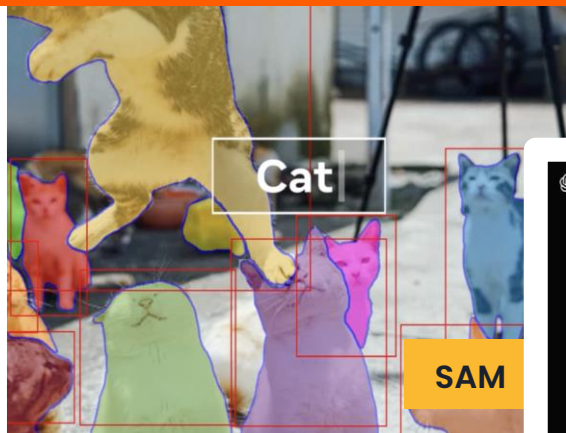🏆 **Kaggle:** Load and analyze any computer vision datasets from Kaggle. Get ahead of your competition with data insights.

🔗 Learn More.

Visual Layer

# From fastdup to Visual-Layer



**Explore**

**Enrich**

**Extract**

Visual Layer

# Explore - Enrich - Extract

**Explore:** Search and discover relevant data blazingly fast.

**Enrich:** Enrich your dataset with state-of-the-art AI models.

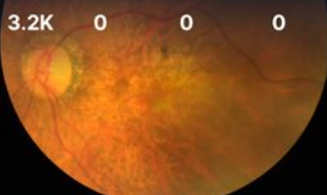**Extract:** Create meaningful subsets for downstream pipelines.

Visual Layer

# Introducing Public VL-Datasets

| retina | ⋮ |
|---|---|
| Created Aug 14 2023 | |

| Images | Objects | Videos | Video Frames |
|---|---|---|---|
| 3.2K | | 0 | 0 |

| ImageNet21K 100K Issues | ⋮ |
|---|---|
| Created Jun 19 2023 | |

| Images | Objects | Videos | Video Frames |
|---|---|---|---|
| 100K | | | |

| KITTI | ⋮ |
|---|---|
| Created Jul 10 2023 | |

| Images | Objects | Videos | Video Frames |
|---|---|---|---|
| 15K | 12K | 0 | 0 |

| DeepFashion | ⋮ |
|---|---|
| Created Jul 10 2023 | |

| Images | Objects | Videos | Video Frames |
|---|---|---|---|
| 289K | 41.3K | | 0 |

| RVL CDIP | ⋮ |
|---|---|
| Created Aug 9 2023 | |

| Images | Objects | Videos | Video Frames |
|---|---|---|---|
| 39.8K | 0 | 0 | 0 |

| COCO | ⋮ |
|---|---|
| Created Jul 13 2023 | |

| Images | Objects | Videos | Video Frames |
|---|---|---|---|
| 164K | 109K | 0 | 0 |

| CelebA | ⋮ |
|---|---|
| Created Jul 10 2023 | |

| Images | Objects | Videos | Video Frames |
|---|---|---|---|
| 203K | 0 | 0 | 0 |

| Oxford-IIIT Pet | ⋮ |
|---|---|
| Created Jul 10 2023 | |

| Images | Objects | Videos | Video Frames |
|---|---|---|---|
| 7.4K | 0 | 0 | 0 |

| ImageNet-1K | ⋮ |
|---|---|
| Created Jun 21 2023 | |

| Images | Objects | Videos | Video Frames |
|---|---|---|---|
| 1.3M | 0 | 0 | 0 |

| Food-101 | ⋮ |
|---|---|
| Created Jul 10 2023 | |

| Images | Objects | Videos | Video Frames |
|---|---|---|---|
| 101K | 0 | 0 | 0 |

# Public VL-Datasets: Coco Mislabels

# Public VL-Datasets: Food-101 Outliers & Batman

# **Public VL-Datasets:** Imagenet 1K semantic search

# Download Quality Issues [Here](#)

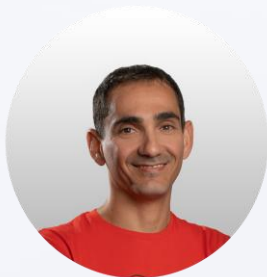| VL Dataset Card | Original Dataset | Explore | Issues CSV | Hugging Face Dataset |
|---|---|---|---|---|
| vl-imagenet-21k | ImageNet-21K | VL PROFILER EXPLORE | DOWNLOAD | 🤗 DATASET WIP |
| vl-imagenet-1k | ImageNet-1K | VL PROFILER EXPLORE | DOWNLOAD | 🤗 DATASET CLICK HERE |
| vl-laion-1b | LAION-1B | VL PROFILER EXPLORE | DOWNLOAD | 🤗 DATASET WIP |
| vl-kitti | KITTI | VL PROFILER EXPLORE | DOWNLOAD | 🤗 DATASET WIP |
| vl-coco | COCO | VL PROFILER EXPLORE | DOWNLOAD | 🤗 DATASET WIP |
| vl-deepfashion | DeepFashion | VL PROFILER EXPLORE | DOWNLOAD | 🤗 DATASET WIP |
| vl-celeba-hq | CelebA-HQ | VL PROFILER EXPLORE | DOWNLOAD | 🤗 DATASET CLICK HERE |
| vl-places365 | Places365 | VL PROFILER EXPLORE | DOWNLOAD | 🤗 DATASET WIP |
| vl-food-101 | Food-101 | VL PROFILER EXPLORE | DOWNLOAD | 🤗 DATASET CLICK HERE |
| vl-oxford-iiit-pet | Oxford-IIIT Pet | VL PROFILER EXPLORE | DOWNLOAD | 🤗 DATASET CLICK HERE |

Visual Layer

# Our **Team**

## Dr. Danny Bickson
### Co-founder & CEO

Co-founder and VP EMEA of Turi (acquired by Apple). CMU Researcher. Sr. Mgr at Apple.

## Dr. Amir Alush
### Co-founder & CTO

Co-founder and CTO of Brodmann17. Highly experienced in building CV/AI Groups and leading into production.

## Prof. Carlos Guestrin
### Co-founder & CSO

Stanford Prof, Co-founder and CEO of Turi (acquired by Apple). Sr. Dir. at Apple. Deep Learning Infra Team Pioneer.

GraphLab | Carnegie Mellon University | turi | 🍎

Brodmann17 | Quris.ai | Bar-Ilan University

dmlc XGBoost | Stanford University | 🍎 | turi | tvm | Carnegie Mellon University

# Our **Team**

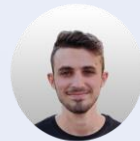**Liah Romantsev**
Technical Customer support

**Sinai Yoktan**
Sr. Backend Engineer

**Gal Bar Nissan**
Staff Engineer

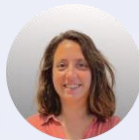**Tom Shani**
Machine Learning Engineer
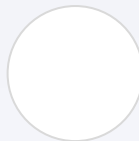
**Dima Frid**
VP Of Engineering

**Elad Yaakov**
Director of Product

**Daniella Bromkish**
Sr. Frontend Engineer
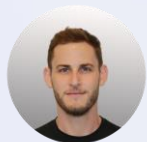
**Head of AI**
WE ARE HIRING!

**Noa Avidar**
Backend Engineer

**Gagandeep Gambhir**
Sr. FrontEnd Engineer

**Guy Singer**
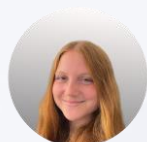Machine Learning Engineer
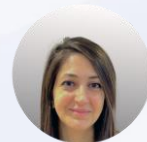
**Achiya Jerbi**
Machine Learning Engineer

**Eli Heifetz**
Principal Engineer

**Ofri Assif**
Office Manager

**Etti Leibovitz**
Head of Product

**Nimrod Mozes**
Design Lead

Backed **By**

Madrona

INSIGHT
PARTNERS

Visual Layer

Thank you

Start exploring today