



KDD2024
BARCELONA, SPAIN



Alibaba



Tongyi

Multi-modal Data Processing for Foundation Models: Practical Guidances and Use Cases



Daoyuan Chen



Yaliang Li



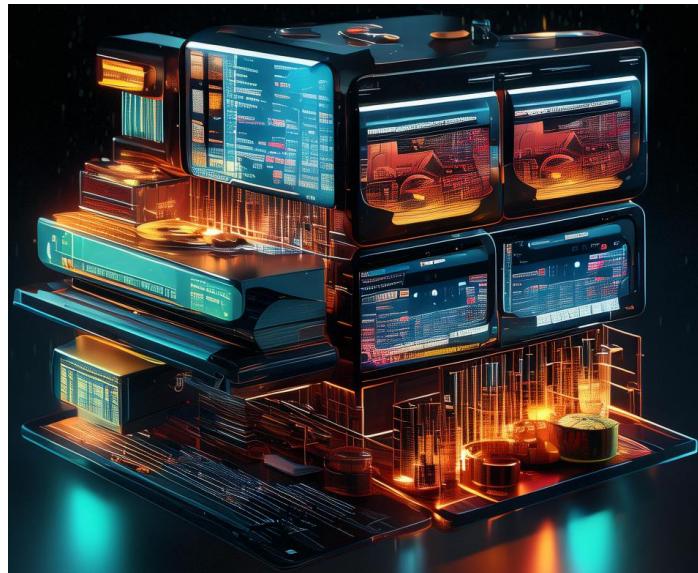
Bolin Ding

Data-Juicer Team

Tongyi Lab, Alibaba Group

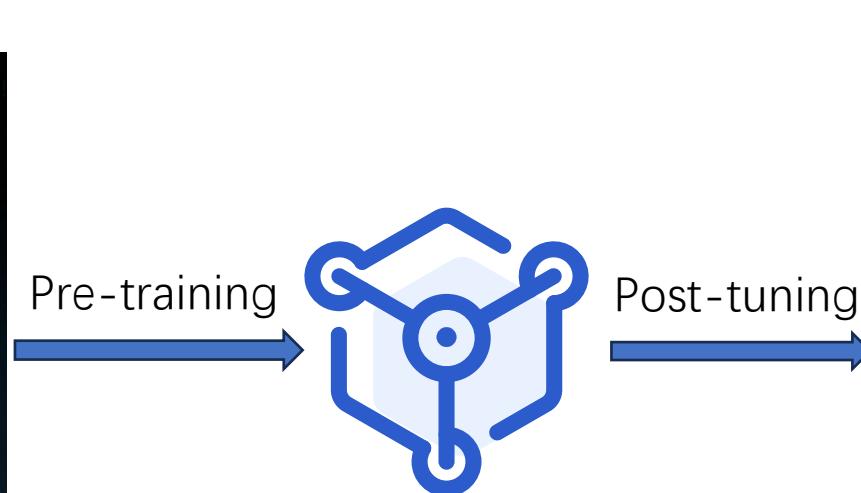
GitHub: <https://github.com/modelscope/data-juicer>

Multi-modal Foundation Model



Multi-modal Data

Video, Image, Text, Voice,



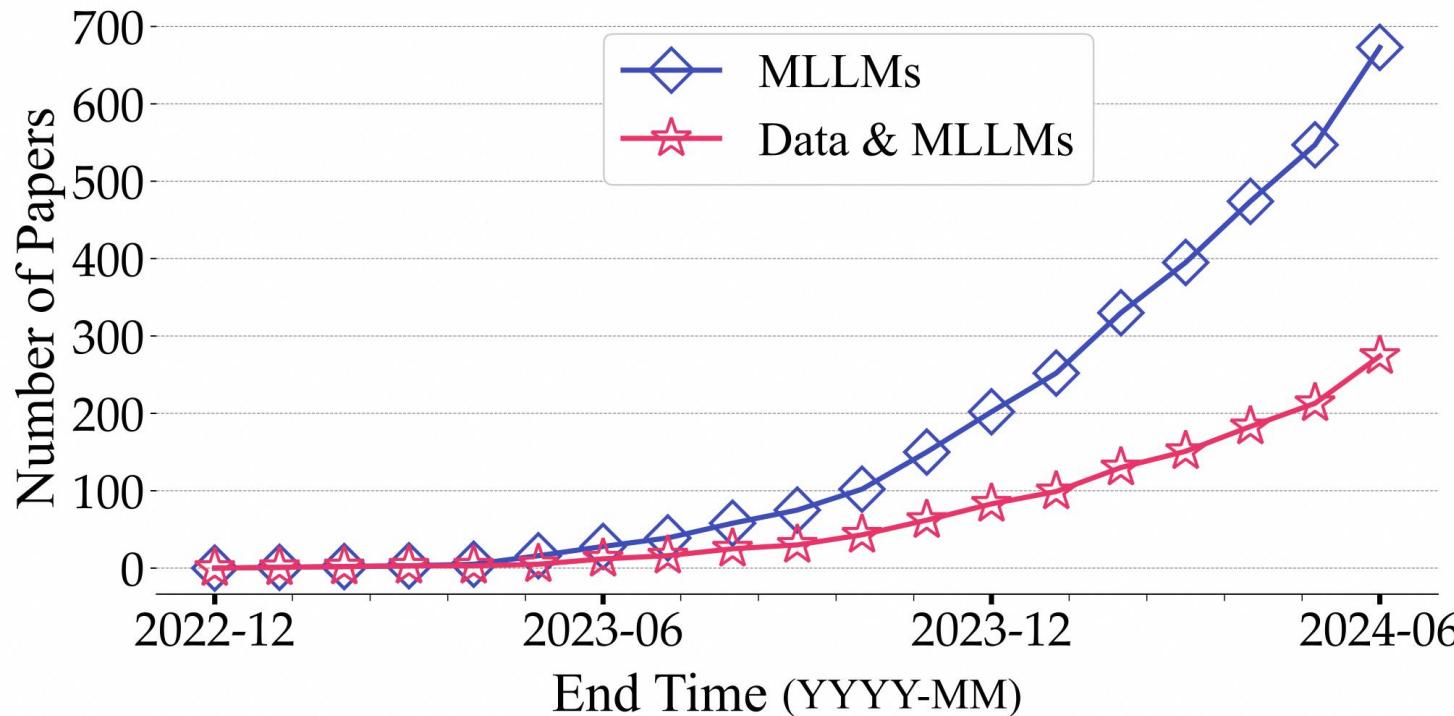
Multi-modal Large Language Models (MLLMs)

GPT4v, Sora, Qwen-VL, Gemini,

Applications

- AI Generative Content
- Multi-modal Dialogue
- Film Editor
- Education
- Entertainment
- Media
-

Multi-modal Foundation Model



Model-centric → Data-centric → Data-Model Co-development

[1] (arXiv:2407.08583) The Synergy between Data and Multi-Modal Large Language Models: A Survey from Co-Development Perspective

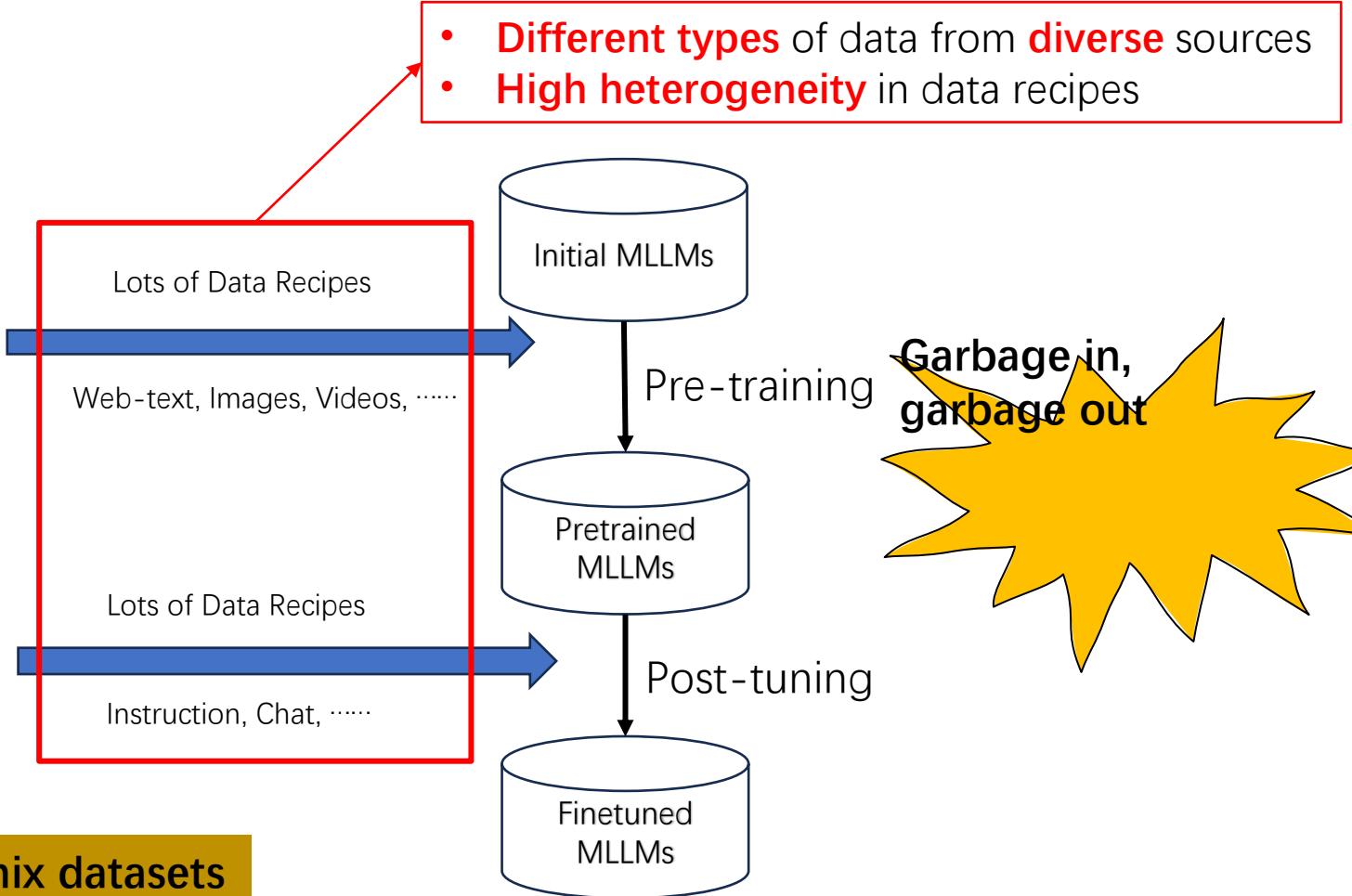
From Data to MLLMs



Hugging Face Search models, datasets, users...
Main Tasks Libraries Languages Licenses Other
Datasets 10,161
Modalities Reset Modalities
3D Audio Geospatial Image
Tabular Text Time-series Video
Size (rows)
<1K >1T

A screenshot of the Hugging Face platform's dataset search interface. It shows a search bar, a main menu with options like Main, Tasks, Libraries, Languages, Licenses, and Other, and a list of datasets. A red box highlights the 'Datasets' section, which shows 10,161 results. Another red box highlights the 'Modalities' section, which includes options for 3D, Audio, Geospatial, Image, Tabular, Text, Time-series, and Video. A slider for 'Size (rows)' is also visible.

Data Recipe: process & mix datasets



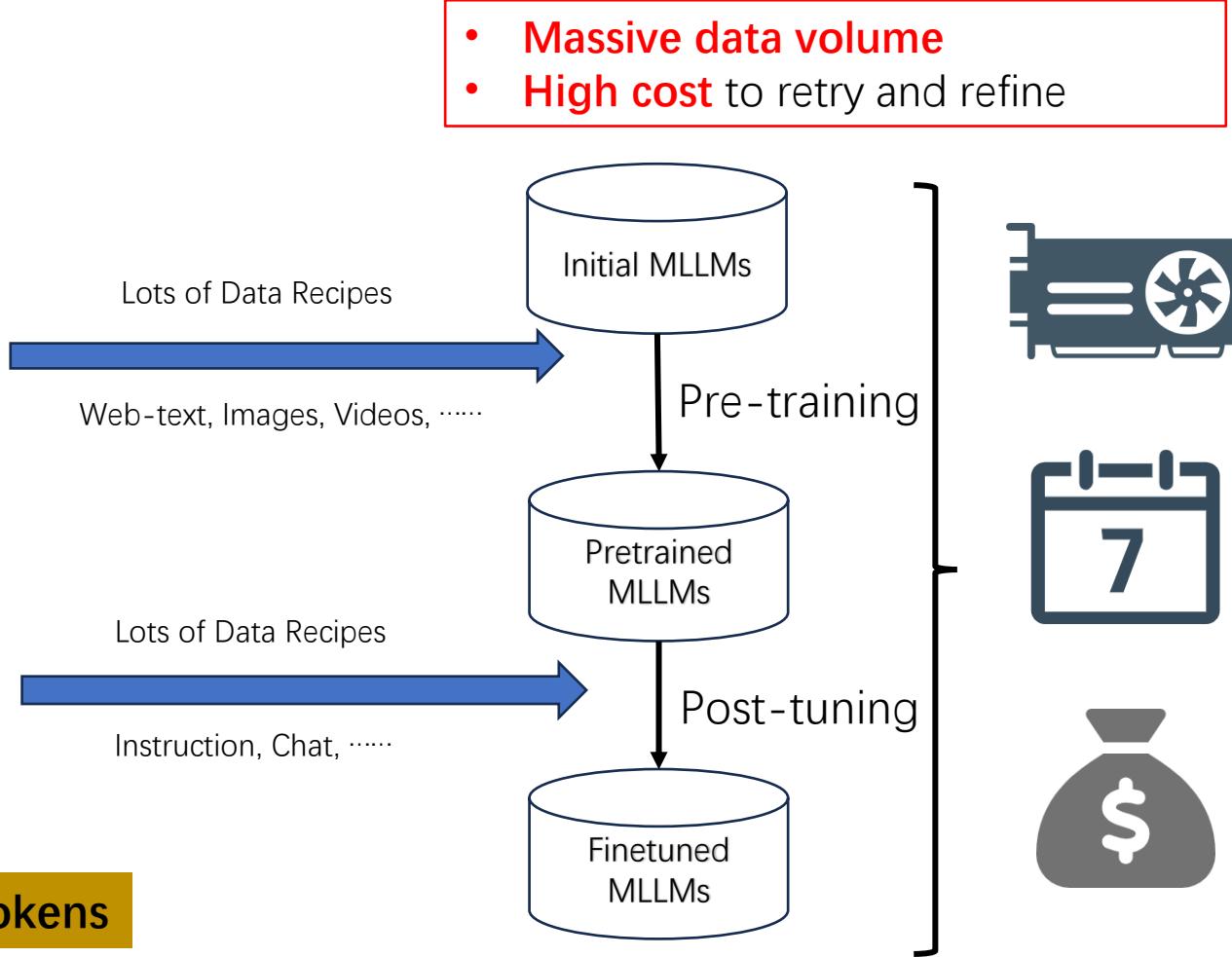
- [2] (SIGMOD'24) Data-juicer: A one-stop data processing system for large language models
[3] (arXiv:2405.14908) Data Mixing Made Efficient: A Bivariate Scaling Law for Language Model Pretraining

From Data to MLLMs



The screenshot shows the Hugging Face Model Hub interface. At the top, there's a search bar with the placeholder "Search models, datasets, users...". Below it, a navigation bar includes "Main", "Tasks", "Libraries", "Languages", "Licenses", and "Other". A red box highlights the "Datasets" section, which displays a count of 10,161 datasets. Another red box highlights the "Modalities" filter, which includes options like "3D", "Audio", "Geospatial", "Image" (selected), "Tabular", "Text" (selected), "Time-series", and "Video". A slider for "Size (rows)" ranges from "<1K" to ">1T". To the right, there are three dataset cards: "lmmslab/LLaVA-", "multimodalart/1", and "jainr3/diffusic", each with a "Viewer" link and an "Updated" timestamp.

Billion/Trillion Tokens



The Data-Juicer System

➤ Systematic & Reusable

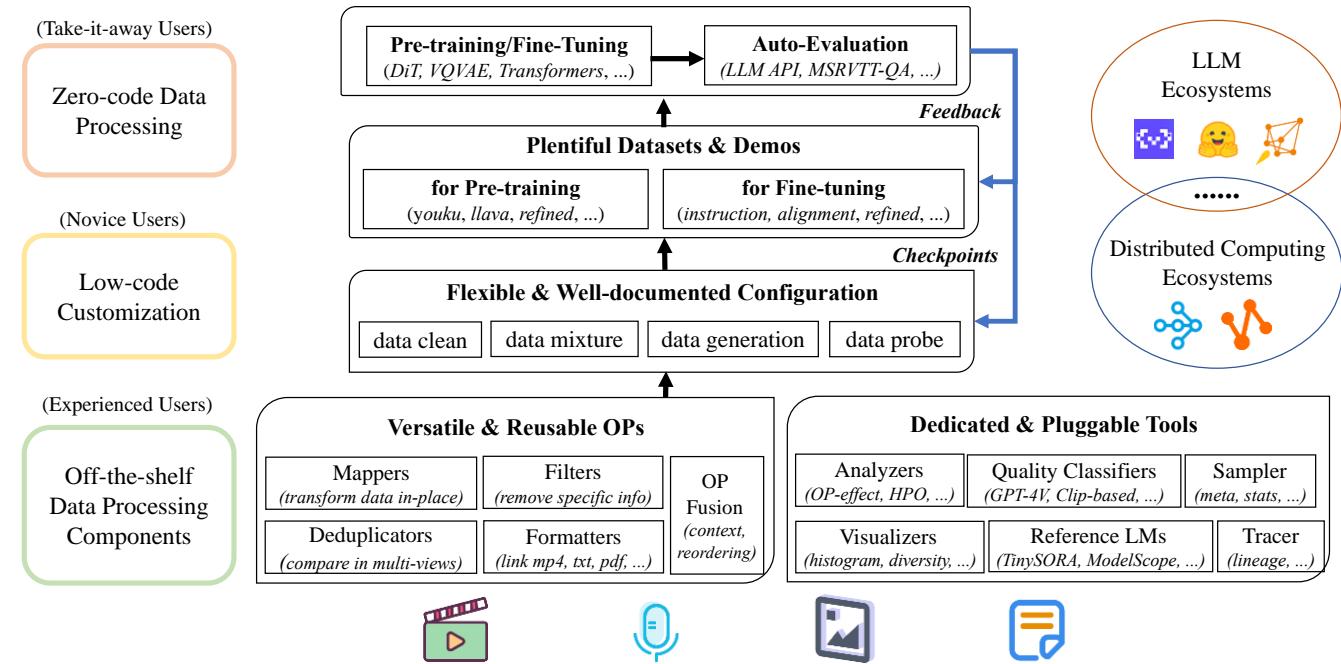
- 100+ standard operators
- 50+ tools & data recipes
- Data analysis, cleaning, synthesis

➤ Efficient & Actionable

- Industrial performance optimization
- Distributed data processing
- Data-model co-development

➤ Effect-Proven & Instructional

- Rank-1 text-to-video model on VBench
- Beats GPT4v/Gemini on MMVP bench
- Maintained survey & practical insights



[2] (SIGMOD'24) Data-juicer: A one-stop data processing system for large language models

Tutorial Outline

- **Data Processing Foundations**
 - Building Blocks of Data Processing: Data-Juicer's Operators
 - Composing Atomic Capabilities: Data-Juicer's Data Recipes
- **Advanced Data Processing**
 - Exploring Data Recipes: The Data-Juicer Sandbox Lab
 - From Exploration to Production: High-Performance Data Factory
- **Use Cases:** From Text to Video Data Processing
- **Resources and Conclusion**

Jupyter
notebooks
embedded

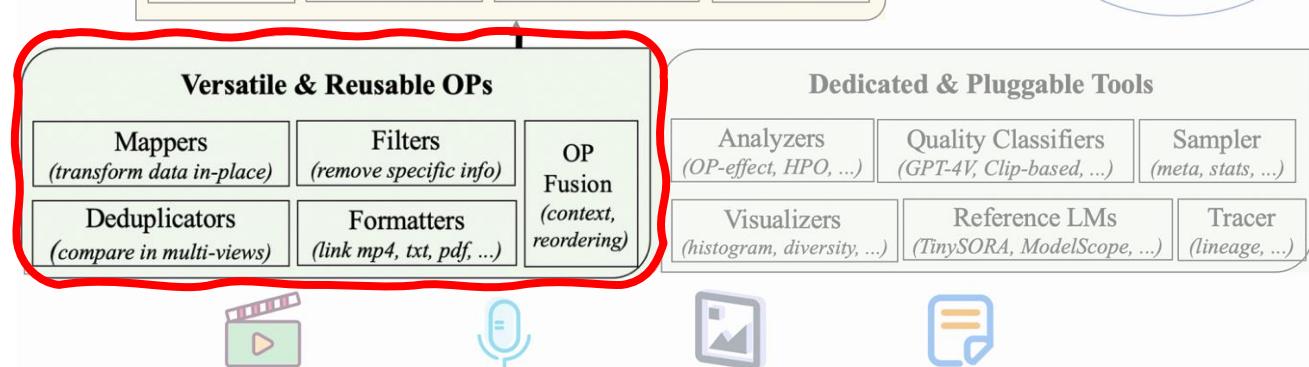
Tutorial Outline

- Data Processing Foundations
 - **Building Blocks of Data Processing: Data-Juicer's Operators**
 - Composing Atomic Capabilities: Data-Juicer's Data Recipes
- Advanced Data Processing
 - Exploring Data Recipes: The Data-Juicer Sandbox Lab
 - From Exploration to Production: High-Performance Data Factory
- Use Cases: From Text to Video Data Processing
- Resources and Conclusion

Data-Juicer: Standard OPs

- **Standard** Interface
- Basic **capabilities** of data processing
- Input: datasets
- Output: processed datasets

Type	Num.	Description
Formatter	7	Discovers, loads, canonicalizes source data
Mapper	46	Edits and transforms samples
Filter	43	Filters out low-quality samples
Deduplicator	5	Detects and removes duplicate samples
Selector	4	Selects top samples based on ranking



From Formatters to Dataset

- Unified format, path-based storage
- Text as anchor, expressing relationships
 - Alignment & Position
 - Local & Global
- Built-in typical dataset conversion tools
 - Video-ChatGPT <--> Data-Juicer
 - Youku-CN <--> Data-Juicer
 - InternVid <--> Data-Juicer
 - LLaVA <--> Data-Juicer
 -

```
1 ▶ {
2   "text": "<|vid> 欢迎来到杭州! <|eoc|> <|vid> 这是美丽的杭州西湖，壮观的钱塘江大潮。 <|eoc|> 让我们目光聚焦到市区西北部。 <|eoc|> <|aud> 聚焦到阿里云! <|vid> <|eoc|> <|img> 这里则是坐落于杭州市西湖区的阿里巴巴云谷园区。 <|eoc|>",
3   "videos": [
4     "path/to/the/video/of/hangzhou",
5     "path/to/the/video/of/WestLake-and-QianTang-River",
6     "path/to/the/video/of/CloudValley"
7   ],
8   "audios": [
9     "path/to/the/audio/of/welcome-to-alibaba-cloud"
10  ],
11  "images": [
12    "path/to/the/image/of/AlibabaCloud-logo",
13  ],
14
15
16  "meta": {
17    "src": "customized",
18    "version": "0.1",
19    "author": "xxx"
20  },
21  "stats": {
22    "lang": "zh",
23    "lang_score": 0.9645169377,
24    "video_frames_text_matching_score": 0.8,
25    "video_duration": 15.5,
26    "avg_audio_image_match_score": 0.7,
27    "image_aesthetics_scores": 0.6
28  }
29 }
```

Data-Juicer: Standard OPs

- Standard OP base class

```
1  class OP:  
2      ...  
3      def __call__(self, dataset: DJDataset) -> DJDataset:  
4          ...  
5          return self.run(dataset, **args)  
  
6  
7  class Mapper(OP):  
8      ...  
9      def process(self, sample: Dict) -> Dict:  
10         # process a single sample  
11         ...  
12  
13         def run(self, dataset: DJDataset) -> DJDataset:  
14             return dataset.map(self.process, **args)  
15         ...  
16  
17  class Filter(OP):  
18      ...  
19      def compute_stats(self, sample: Dict) -> Dict:  
20          # compute stats of a single sample  
21          ...  
22  
23      def process(self, sample: Dict) -> bool:  
24          # decide whether to keep this single sample  
25          ...  
26  
27      def run(self, dataset: DJDataset) -> DJDataset:  
28          return dataset.map(self.compute_stats, **args)  
29              .filter(self.process, **args)  
30          ...  
31
```

- Flexible invoking logic

```
1  op1_config = ...  
2  op2_config = ...  
3  
4  dataset = Dataset.from_jsonl(path=..., mode="standalone/ray")  
5  
6  op1 = Operator_1(**op1_config)  
7  op2 = Operator_2(**op2_config)  
8  
9  # single OP  
10 dataset = dataset.process(op1)  
11 dataset = dataset.process(op2)  
12  
13 # multiple OPs  
14 dataset = dataset.process(op1).process(op2)  
15 # or  
16 dataset = dataset.process([op1, op2])  
17  
18 dataset.export(path=..., format="jsonl")
```

Data Processing Hands-on

- Jupyter notebooks part 1
 - [1.1 OP insights](#) about specific OP types, and multi-modal OP examples
 - [1.2 datasets loading](#): multi-modal dataset format <--> Data-Juicer format
 - [1.3 Brief introduction](#) of how to develop a new OP

Tutorial Outline

- **Foundational Abilities**
 - Building Blocks of Data Processing: Data-Juicer's Operators
 - **Composing Atomic Capabilities: Data-Juicer's Data Recipes**
- Advanced Data Processing
 - Exploring Data Recipes: The Data-Juicer Sandbox Lab
 - From Exploration to Production: High-Performance Data Factory
- Use Cases: From Text to Video Data Processing
- Resources and Conclusion

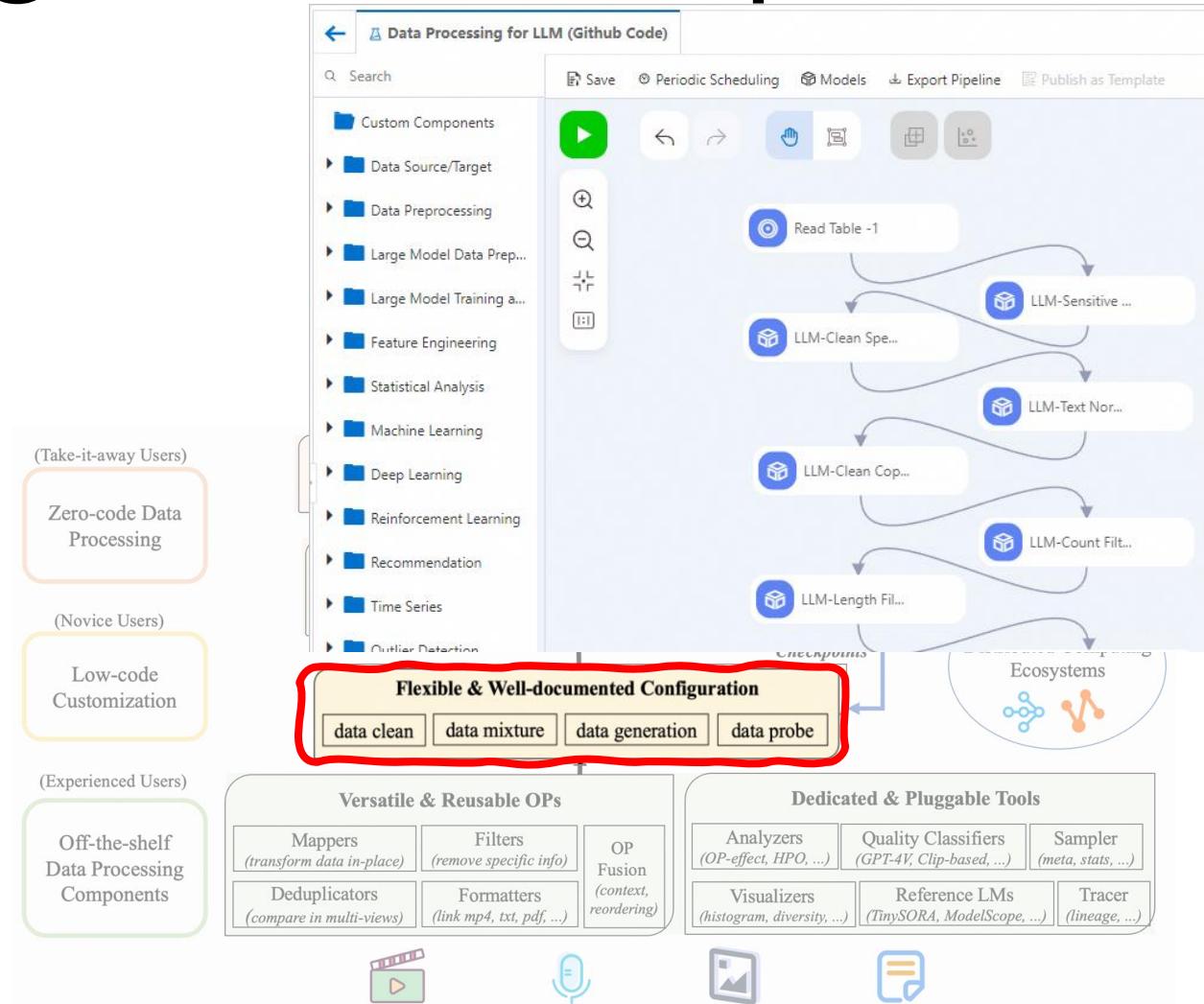
Data-Juicer: Data Recipes

- Data recipe includes:
 - Data source
 - Mixture ratios
 - Data processing pipeline
 -
- **Effect-proven**
- **Reusable** and **shared** for the community



How to config for data recipes

- Three ways to config:
 - A piece of code: describe the processing pipeline
 - A config file (**recommended: simple & flexible**)
 - Drag-drop interface: equivalent to a config file, (available on PAI-Designer)



Example: A Config File for Data Recipe

```
1 project_name: 'video-dj-recipe-demo'
2 dataset_path: 'video-dj-recipe-youku.jsonl'
3 export_path: 'video-dj-recipe-youku-generated.jsonl'
4 np: 42
5 open_tracer: true
6 video_key: 'videos'
7 video_special_token: '<video>'
8 eoc_special_token: '<|__dj__eoc|>'

9
10 # process schedule: a list of several process operators with their arguments
11 process:
12   - video_deduplicator:
13   - fix_unicode_mapper:
14   - punctuation_normalization_mapper:
15
16   - video_split_by_scene_mapper:
17     detector: 'ContentDetector'
18     threshold: 27.0
19     min_scene_len: 15
20     show_progress: false
21   - video_captioning_from_audio_mapper:
22     keep_original_sample: true
23   - video_tagging_from_audio_mapper:
24     hf_dst: 'MIT/ast-finetuned-audio-set-10-10-0.4573'
25   - video_frames_text_similarity_filter:
26     hf_clip: 'openai/clip-vit-base-patch32'
27     min_score: 0.1
28     max_score: 1.0
29     frame_sampling_method: 'all_keyframes'
30     frame_num: 3
31     horizontal_flip: false

32
33   # number of subprocess to process your data
34   # Key name of field to store the list of samples
35   # The special token that represents an video
36   # The special token that represents the end of a video

37
38   # deduplicator to deduplicate samples at document level
39   # fix unicode errors in text.
40   # normalize unicode punctuations to English

41   # split videos into scene clips
42   # PySceneDetect scene detector. Should be one of the detectors
43   # threshold passed to the detector
44   # minimum length of any scene
45   # whether to show progress from scenedetect
46   # caption a video according to its audio streams
47   # whether to keep the original sample. If it's False, the sample will be discarded.
48   # Mapper to generate video tags from audio stream
49   # Huggingface model name for the audio classifier
50   # keep samples those similarities between sampled frames are above the threshold
51   # clip model name on huggingface to compute the similarity
52   # the min similarity to keep samples.
53   # the max similarity to keep samples.
54   # sampling method of extracting frame images from a video
55   # the number of frames to be extracted uniformly
56   # flip frame image horizontally (left to right)
```

Global Arguments:

- Paths to: input & output datasets,
- Num of workers: process in parallel
-

Operators in Ordered List:

- OP name
- Arguments of OPs
- Including several categories of OPs

Filtering (single-modality quality):

- video de-dup; text de-noise;

Enhancing (diversity)

- video clips; re-captioning;

Filtering (cross-modality quality):

- text-video matching degree;

Built-in Data Recipes

dataset	#samples before	#samples after	keep ratio	data recipe link	data link	source
arXiv	1,724,497	1,655,259	95.99%	redpajama-arxiv-refine.yaml	Aliyun ModelScope HuggingFace	Redpajama
..... (Textual Data Recipes)						
MGM pretrain (1.2M)	1,266,268	159,288 (repeat 4 times)	12.58%	mgm-pretrain-top2-ops-refine.yaml	Aliyun ModelScope HuggingFace	LLaVA-1.5 ALLaVA-4V
Video-text	1,217,346	147,176	12.09%	2_multi_op_pipeline.yaml	Aliyun ModelScope HuggingFace	Panda-70M InternVid MSR-VTT
..... (Multi-modal Data Recipes)						

➤ Refined for better diversity

- Combination of SOTA dataset, e.g., Panda-70M, InternVid, MSR-VTT, ...
- Contrastive data synthesis
-

➤ Refined for higher quality

- Data cleaning guided by OPs, tools, sandbox, ...
- Duplicating good data
-

Effect-Proven Data Recipes

dataset	#samples before	#samples after	keep ratio	data recipe link	data link	source
arXiv	1,724,497	1,655,259	95.99%	redpajama-arxiv-refine.yaml	Aliyun ModelScope HuggingFace	Redpajama
..... (Textual Data Recipes)						
MGM pretrain (1.2M)	1,266,268	159,288 (repeat 4 times)	12.58%	mgm-pretrain-top2-ops-refine.yaml	Aliyun ModelScope HuggingFace	LLaVA-1.5 ALLaVA-4V
Video-text	1,217,346	147,176	12.09%	2_multi_op_pipeline.yaml	Aliyun ModelScope HuggingFace	Panda-70M InternVid MSR-VTT
..... (Multi-modal Data Recipes)						

➤ Video-text recipe example^[4]

- Rank-1 model (ours) on Vbench leaderboard

➤ Image-text recipe example^[5]

- 50.7 (ours) v.s. 38.7 (GPT-4V) on MMVP benchmark

[4] (arXiv:2407.11784) Data-Juicer Sandbox: A Comprehensive Suite for Multimodal Data-Model Co-development

[5] (arXiv:2408.04594) ImgDiff: Contrastive Data Synthesis for Vision Large Language Models

Data Processing Hands-on

- Jupyter notebooks part 2
 - [2.1 conduct data processing with existing recipes](#)
 - [2.2 introduction of some basic global information](#)
 - [2.3 how to modify your own recipes](#)

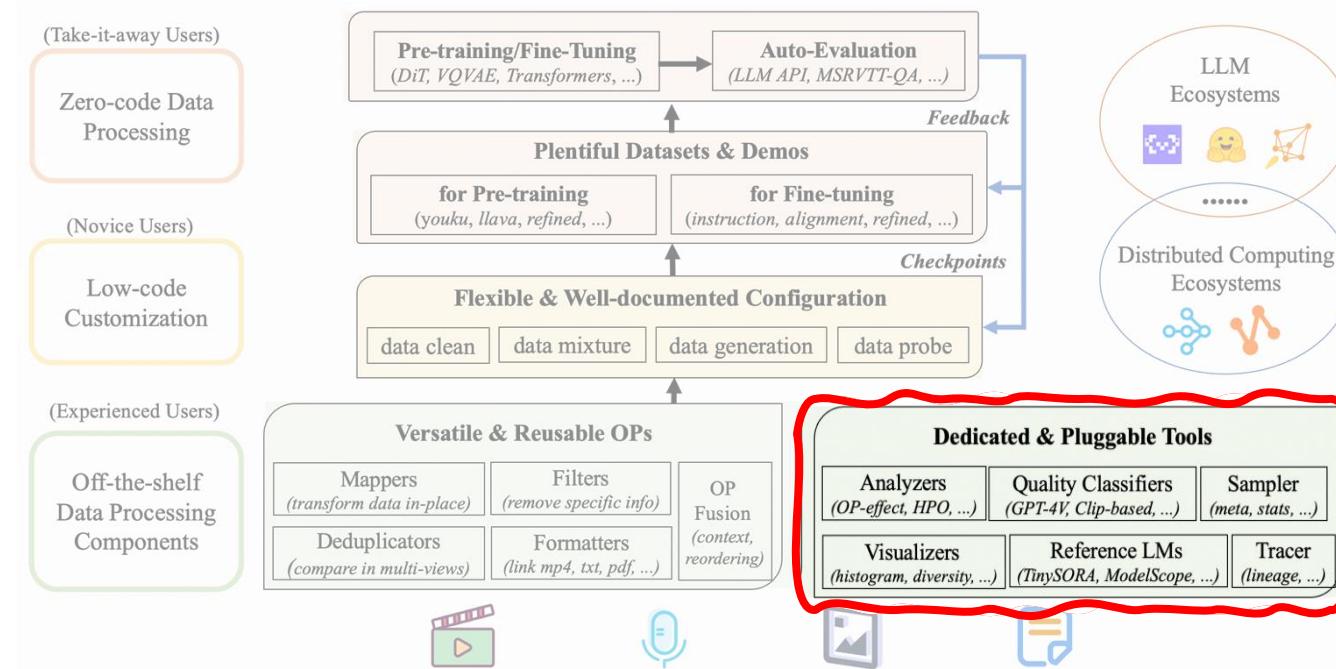
Tutorial Outline

- **Foundational Abilities**
 - Building Blocks of Data Processing: Data-Juicer's Operators
 - Composing Atomic Capabilities: Data-Juicer's Data Recipes
- **Advanced Data Processing**
 - **Exploring Data Recipes: The Data-Juicer Sandbox Lab**
 - From Exploration to Production: High-Performance Data Factory
- **Use Cases:** From Text to Video Data Processing
- **Resources and Conclusion**

Data-Juicer: Tools

➤ **Dedicated** for **complex** and **specific** functions, supporting data recipe exploration

- Analysis
- Processing trace
- Visualization
-



Tools: Tracer

- Sample-level changes before/after each OP

Filter: filtered-out samples

```
⊕{
    "text": "你好, 请问你是谁",
    "meta": ⊕{
        "src": "customized",
        "date": null,
        "version": null,
        "author": "xxx"
    },
    "stats": ⊕{
        "lang": "zh",
        "lang_score": 0.9450979829
    }
}
```

Mapper: modified samples

```
⊕{
    "original_text": "This is the email of our company: euqdh@cjqi.com",
    "processed_text": "This is the email of our company: "
}
```

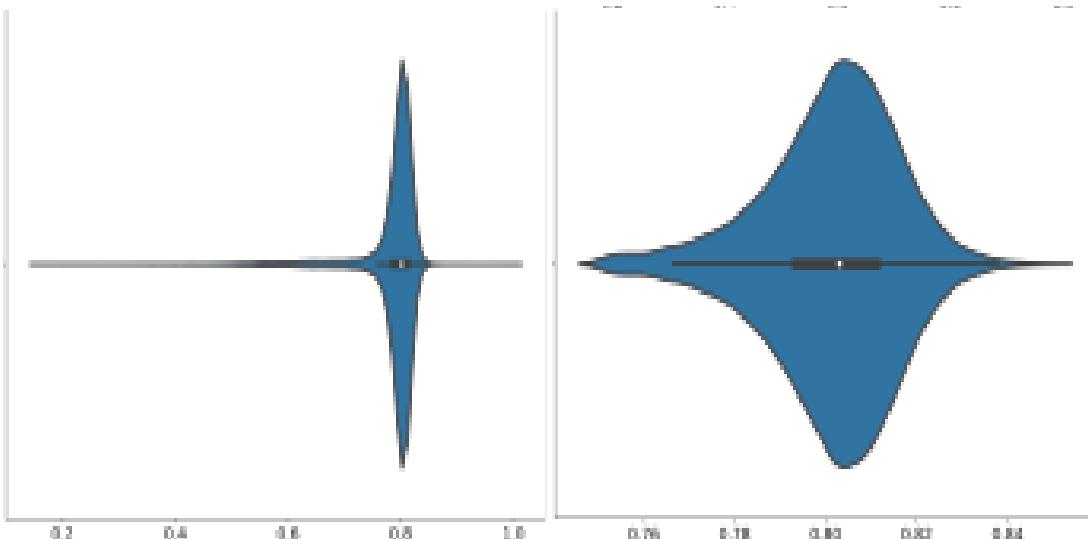
Deduplicator: (near-) duplicate sample pairs

```
⊕{
    "dup1": ⊕{
        "text": "This paper proposed a novel method on LLM pretraining.",
        "meta": ⊕{
            "meta_inner": ⊕{
                "src": "customized",
                "date": null,
                "version": null,
                "author": "xxx"
            }
        },
        "stats": ⊕{
            "lang": "en",
            "lang_score": 0.9568607211
        },
        "hash": "df544ffbc314a6d27b2847429246be76"
    },
    "dup2": ⊕{
        "text": "This paper proposed a novel method on LLM pretraining.",
        "meta": ⊕{
            "meta_inner": ⊕{
                "src": "customized",
                "date": null,
                "version": null,
                "author": "xxx"
            }
        },
        "stats": ⊕{
            "lang": "en",
            "lang_score": 0.9568607211
        },
        "hash": "df544ffbc314a6d27b2847429246be76"
    }
}
```

Tools: Analyzer

➤ **Quality:** statistics distribution

- Character repetition ratio
- Perplexity
-



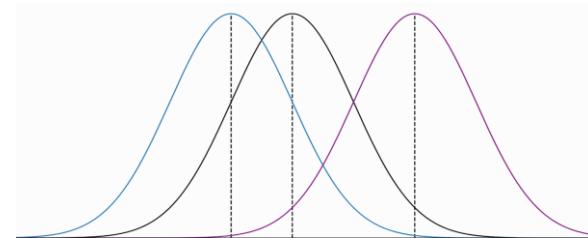
➤ **Diversity:** lexical combination



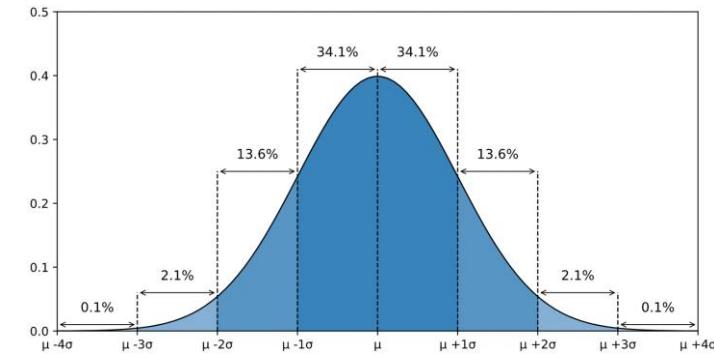
Tools: Recipe Iteration

K-Sigma & Model Sensitivity

- **Analyze** data automatically
 - Compute quantitative values of each feature dimension
 - Statistics (fluency, aesthetic, cross-modal matching, etc.)
 - Sample data and evaluate the sensitivities from models
- **Filter out** data outside the k-sigma range



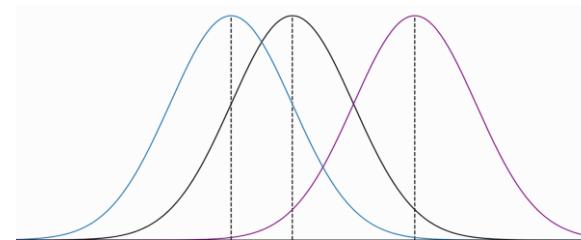
stats from 40+ Data-Juicer Filters



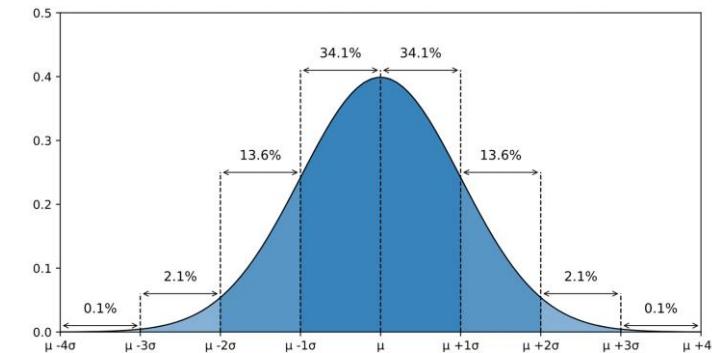
Tools: Recipe Iteration

K-Sigma & Model Sensitivity

- **Analyze** data automatically
 - Compute quantitative values of each feature dimension
 - Statistics (fluency, aesthetic, cross-modal matching, etc.)
 - Sample data and evaluate the sensitivities from models
- **Filter out** data outside the k-sigma range

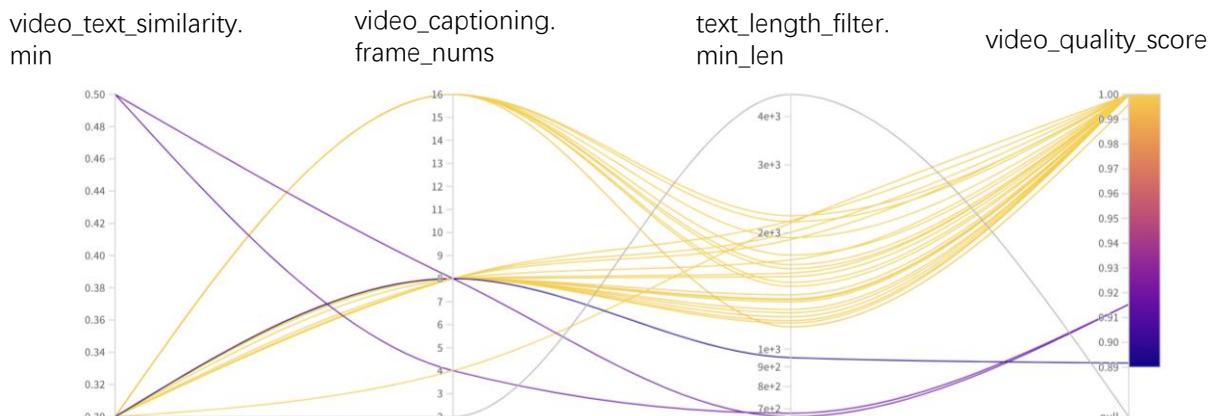


stats from 40+ Data-Juicer Filters



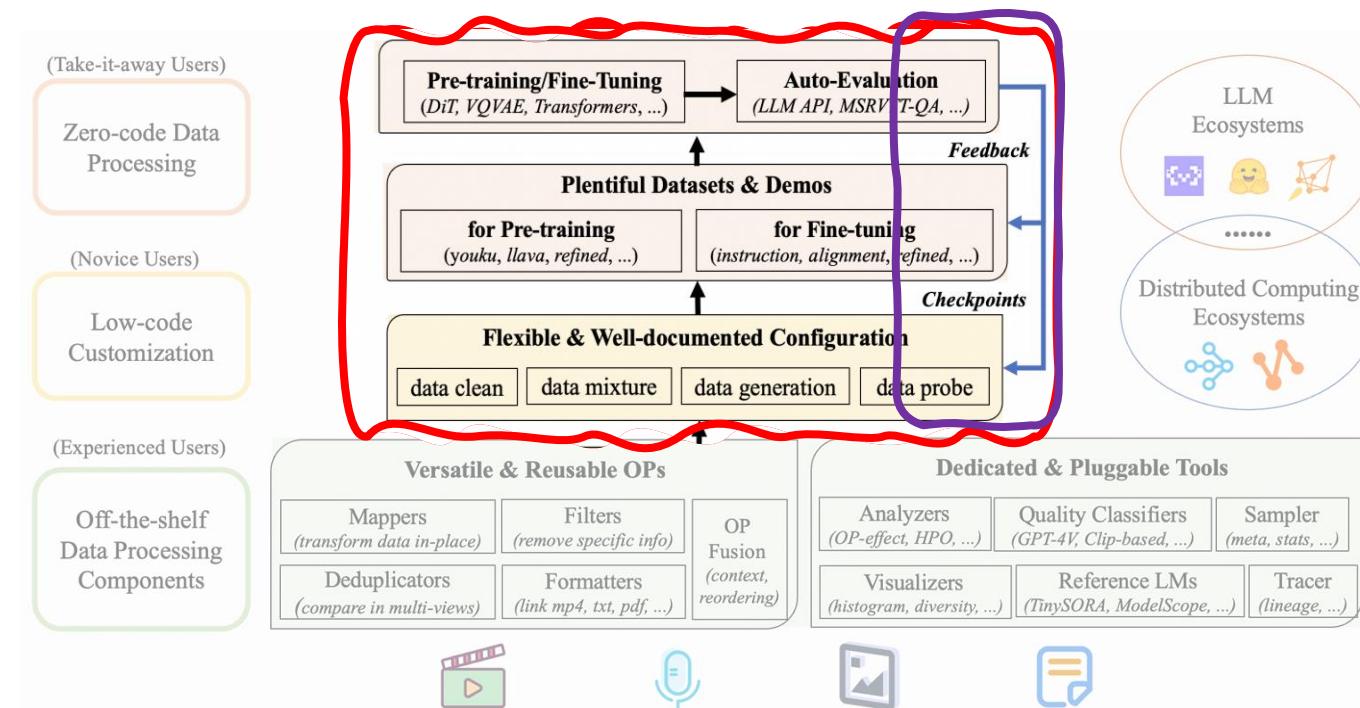
Auto Recipe Optimization

- **Analyze** importance & relevance of hyper-parameters in data recipe
- Integrated W&B Sweeps: support heuristic search & Bayesian opt.



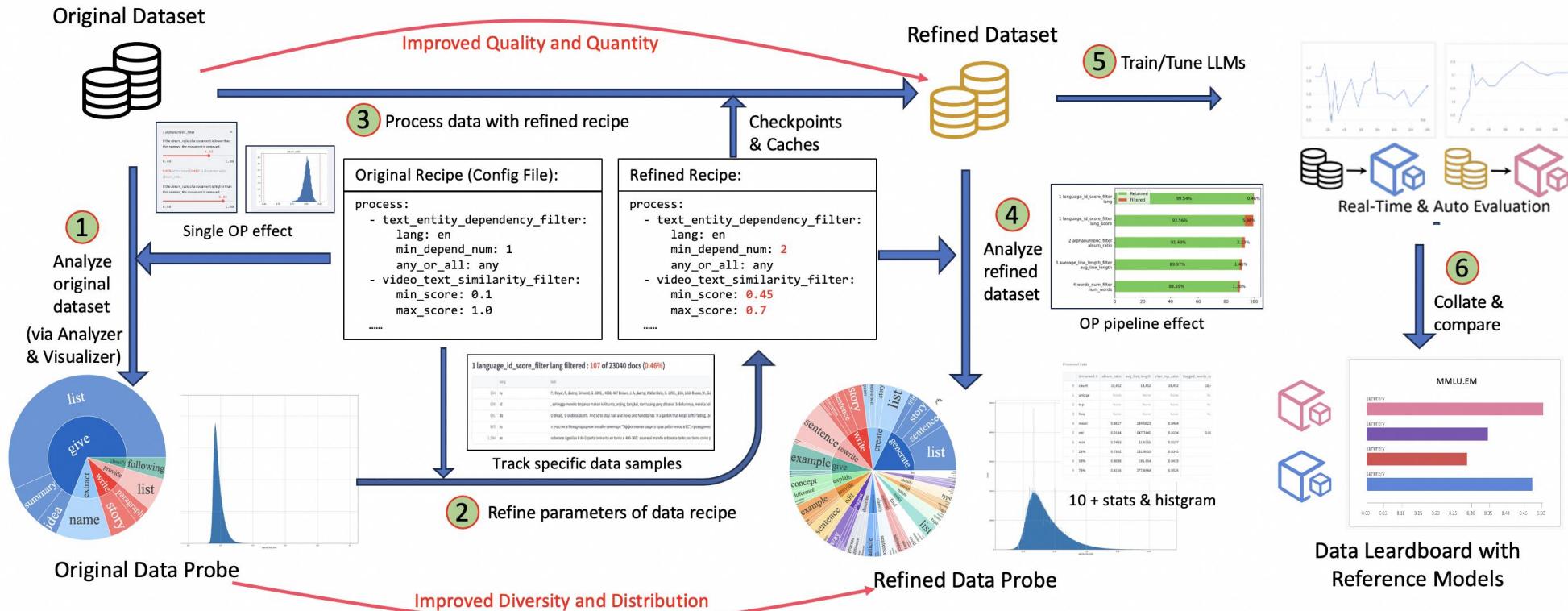
Data-Juicer: Feedback Loop

- Return multiple feedback signals
- Refine recipe systematically



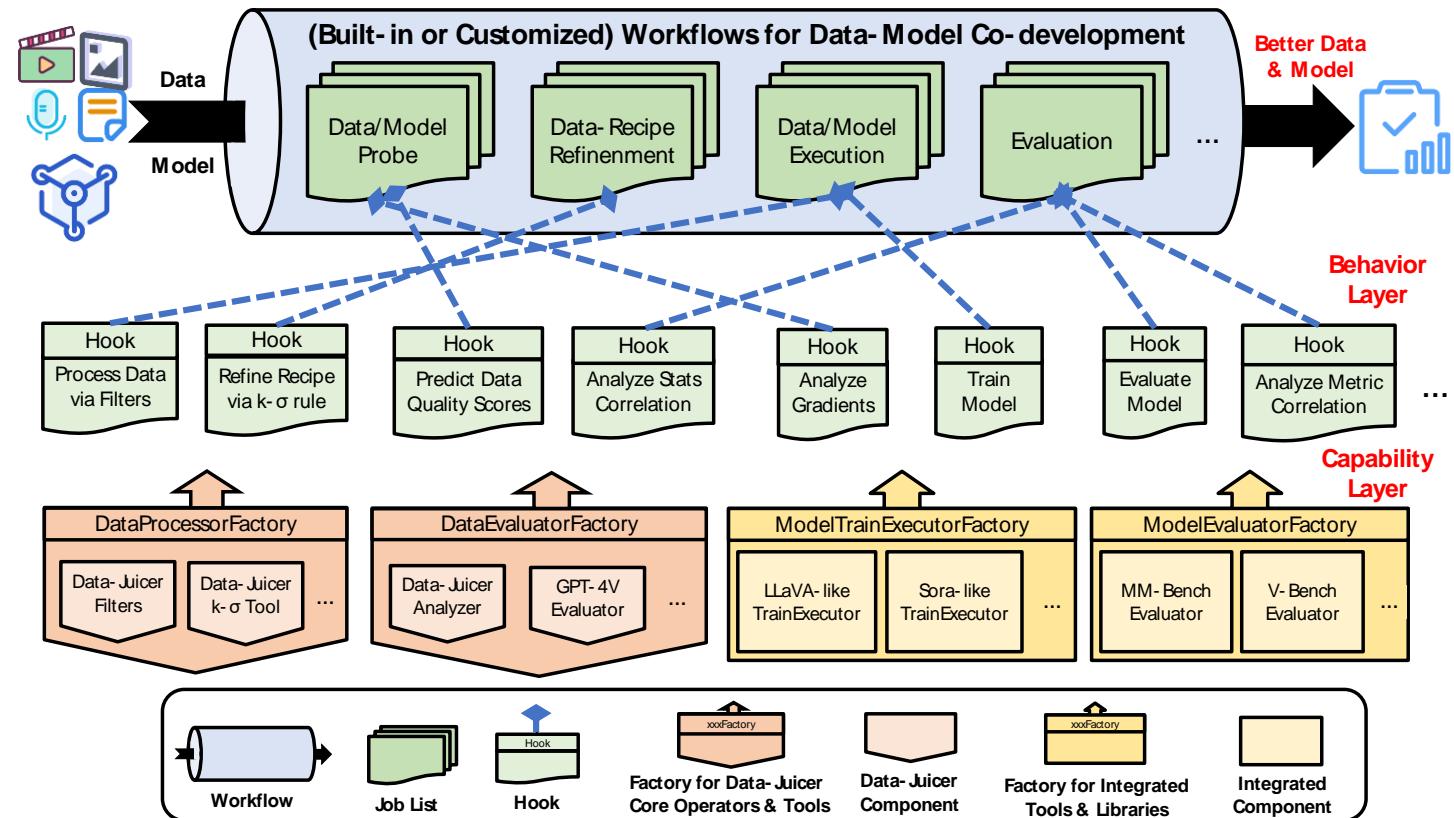
Data-Juicer: Feedback Loop

- Improve data quality with feedbacks from different steps of data-model co-development
- The whole loop has been formed as a **Data-Model Sandbox** in Data-Juicer now!



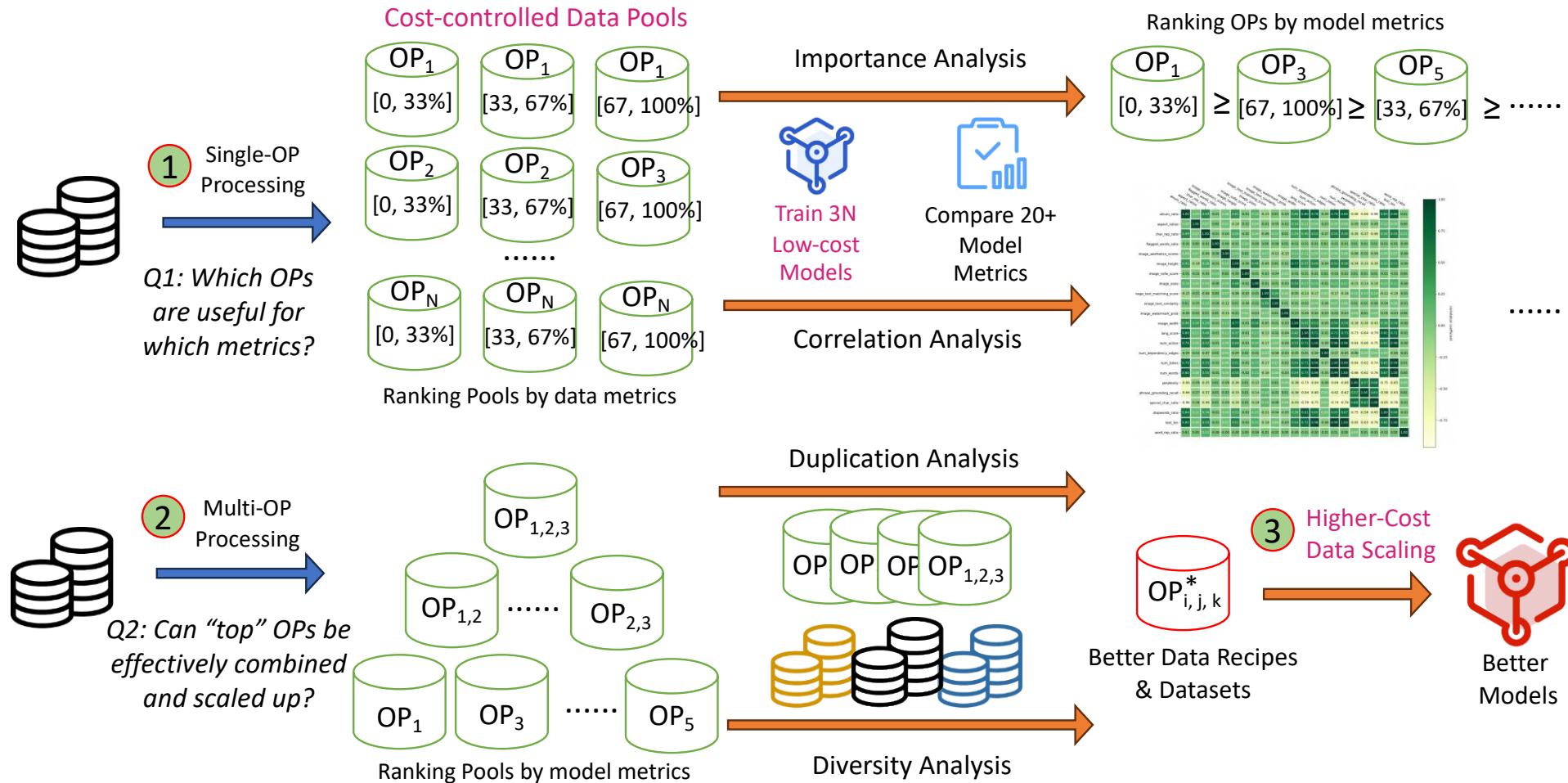
Data-Juicer: Sandbox Lab

- Three-layer Orchestration
 - Workflow
 - Behavior
 - Capability
- Cost-controlled
- Easy Customization



[3] (arXiv:2407.11784) Data-Juicer Sandbox: A Comprehensive Suite for Multimodal Data-Model Co-development

Data-Juicer: Sandbox Lab



Recipe Exploration Hands-on

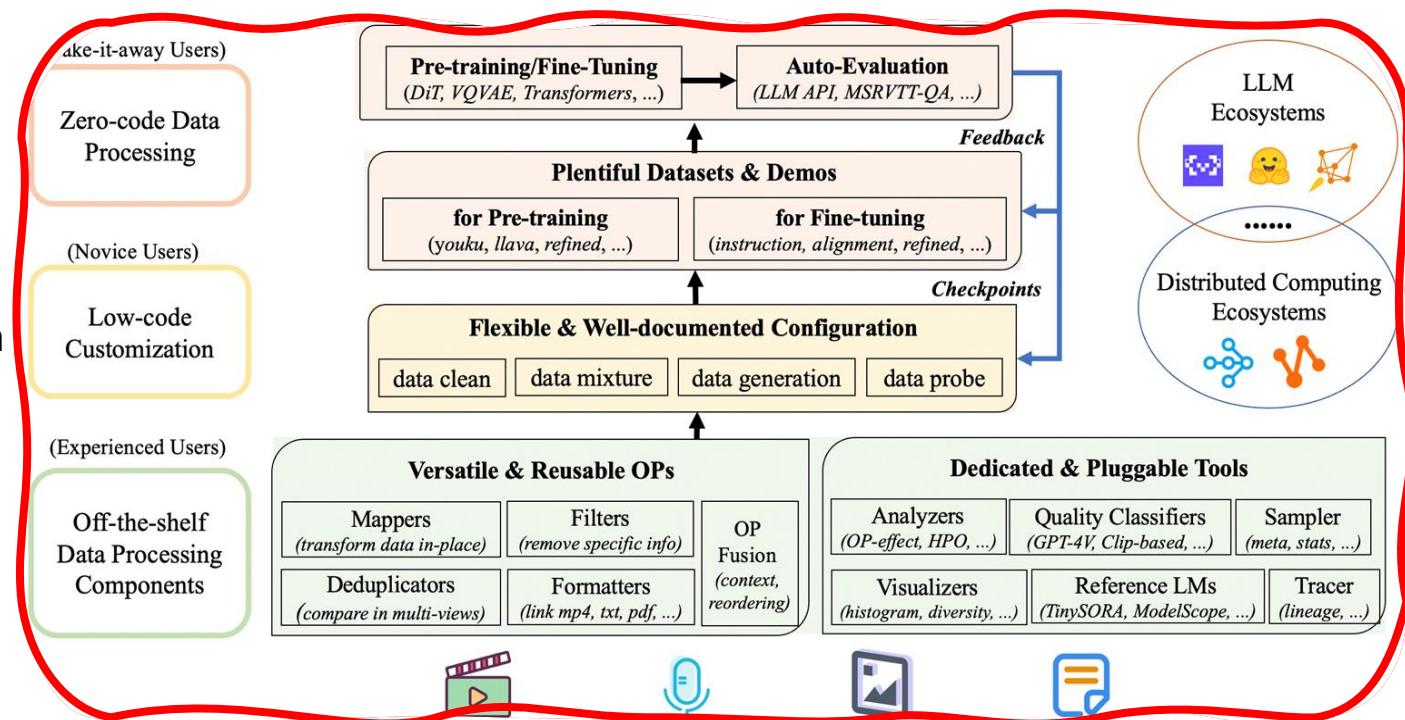
- Jupyter notebooks part 3
 - [3.1 how to use the Data-Juicer tools](#)
 - [3.2 how to use the Data-Juicer SandBox](#)

Tutorial Outline

- **Foundational Abilities**
 - Building Blocks of Data Processing: Data-Juicer's Operators
 - Composing Atomic Capabilities: Data-Juicer's Data Recipes
- **Advanced Data Processing**
 - Exploring Data Recipes: The Data-Juicer Sandbox Lab
 - **From Exploration to Production: High-Performance Data Factory**
- **Use Cases:** From Text to Video Data Processing
- **Resources and Conclusion**

Data-Juicer: Processing Optimization

- **Distributed Processing**
 - Systematic scalability
- **OP Fusion**
 - Avoid redundant computation
- **Resource Management**
 - Adaptively assign resources
- **Fault Tolerance**
 - Increase reliability

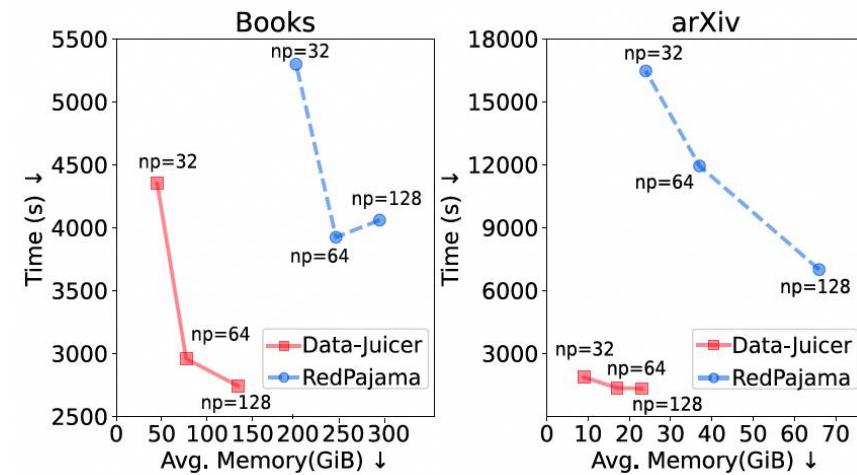
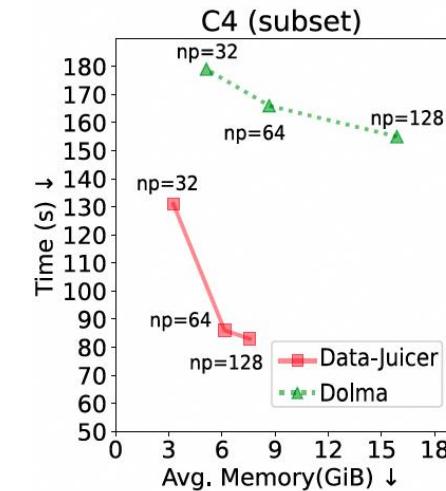


Towards a Production System

Processing Performance

Performance Comparison

- Standalone mode, optimization flags off
- Data Processing Baseline: Dolma & RedPajama
 - Avg. Time ($\downarrow 50.6\%$)
 - Avg. Memory ($\downarrow 55.1\%$)
- The basic implementation is already efficient!

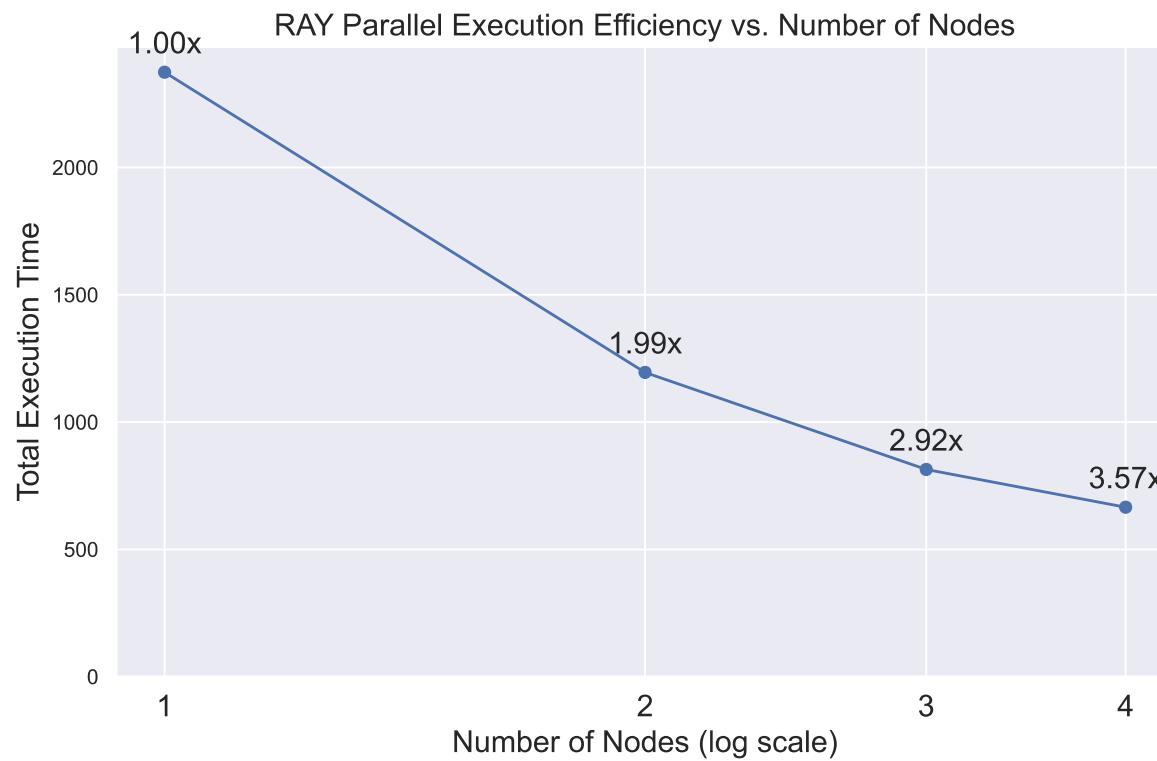


Optimization: Parallelization

Opt 1: Parallel Processing

- Multi-core multi-node
- Support CUDA
- Systematic scalability
with Ray

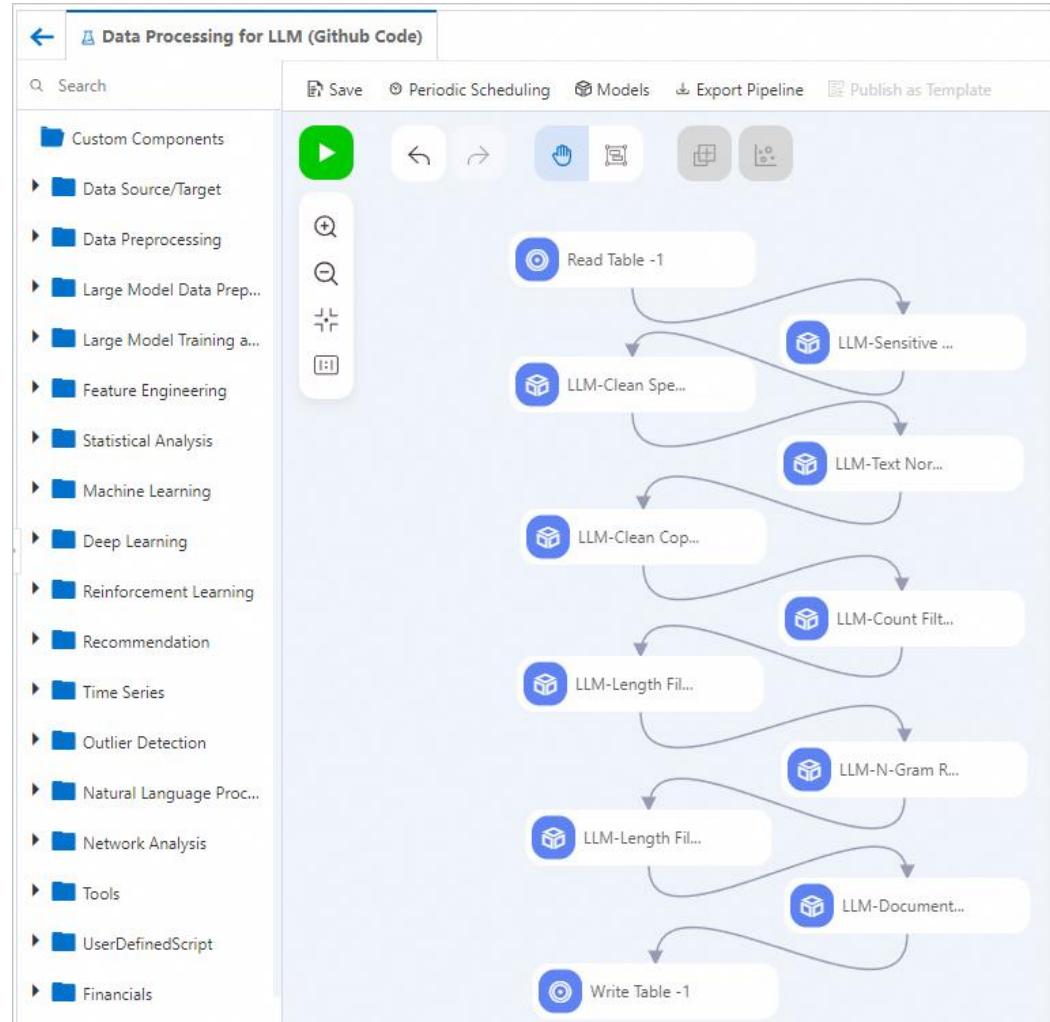
(<https://github.com/ray-project/ray>)



Optimization: Parallelization

Opt 1: Parallel Processing

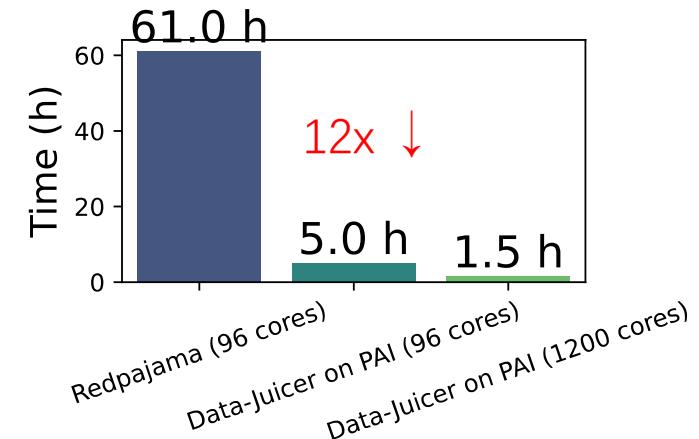
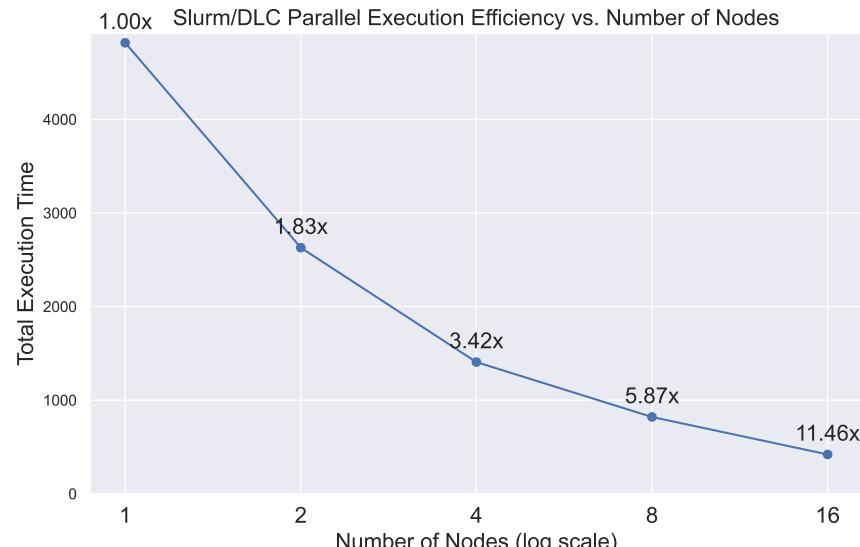
- On-cloud distributed processing:
 - Alibaba Cloud PAI-Designer



Optimization: Parallelization

Opt 1: Parallel Processing

- On-cloud distributed processing:
 - Alibaba Cloud PAI-Designer
 - Alibaba Cloud PAI-DLC/DSW

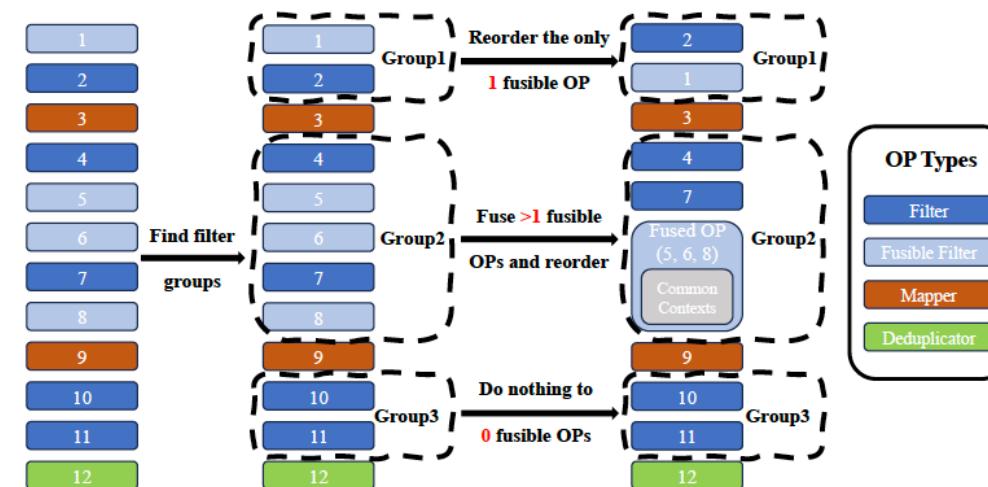


Optimization: OP Fusion

Opt 2: OP Fusion + Context Management + OP Reordering

➤ Core facts:

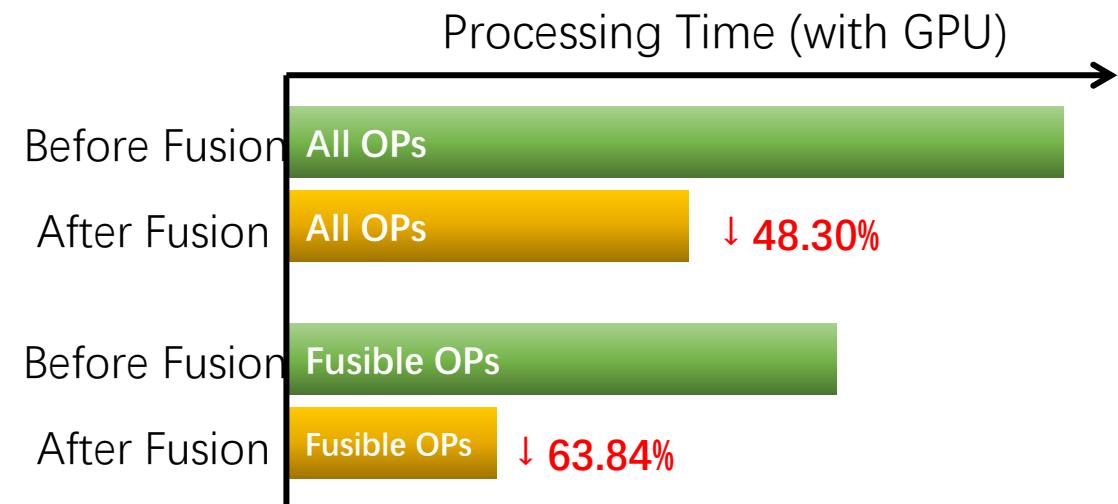
- Some OPs share the same contexts & computation sub-procedures
- Filters are commutative
- Some OPs are much more time-consuming



Optimization: OP Fusion

Opt 2: OP Fusion + Context Management + OP Reordering

- End-to-end Performance Comparsion
- MSR-VTT, 7 OPs (5 capable of fusion)
 - Video Frame-Text Similarity Filter
 - Video NSFW Filter
 - Video Motion Score Filter
 - Video OCR Area Ratio Filter
 - Video Watermark Filter



Optimization: Resource Management

Opt 3: OP-aware parallelization & Quantization

- OP-aware Parallelization
 - different OP requires different amount of resources
 - adaptively adjust parallelization strategy for different OPs
- Low Precision Supported → Less Memory
 - speeding up model-based OPs (with larger batch size)
 - enabling better cleaning/synthesis (with larger model size)

Optimization: Fault Tolerance

Opt 4: Batched Processing & Skipping

- Automatically skip problematic samples → **Worst-case guarantee**
 - Support batched processing for all OPs
 - Dynamic number of returned processed samples

```
2024-07-04 14:40:42 | INFO  | data_juicer.format.mixture_formatter:137 - sampled 4 from 4
2024-07-04 14:40:42 | INFO  | data_juicer.format.mixture_formatter:143 - There are 4 in final dataset
2024-07-04 14:40:42 | INFO  | data_juicer.core.executor:156 - Preparing process operators...
2024-07-04 14:40:42 | INFO  | data_juicer.core.executor:163 - Processing data...
2024-07-04 14:40:42 | WARNING | data_juicer.utils.process_utils:64 - The required CPU number:1 and memory:0GB might be more than the available CPU:12 and memory :6.366420745849609GB. This Op [video_split_by_duration_mapper] might require more resource to run.
video_split_by_duration_mapper_process: 75% | 3/4 [00:09<00:03, 3.69s/ examples]
2024-07-04 14:40:52 | ERROR  | data_juicer.ops.base_op:55 - An error occurred in mapper operation when processing samples {'videos': ['/Users/null/Desktop/Worksapce/Codes/data-juicer/demos/process_video_on_ray/data./videos/video4.mp4'], 'text': ['<_dj_video> 46s videos <|_dj_eoc|>']}, <class 'FileNotFoundException'>: Video [/Users/null/Desktop/Worksapce/Codes/data-juicer/demos/process_video_on_ray/data./videos/video4.mp4] does not exist!
video_split_by_duration_mapper_process: 0% | 0/6 [00:00<?, ? examples/s]
2024-07-04 14:40:52 | INFO   | data_juicer.core.data:162 - OP [video_split_by_duration_mapper] Done in 9.878(s). Left 6 samples.
2024-07-04 14:40:52 | WARNING | data_juicer.utils.process_utils:64 - The required CPU number:1 and memory:0GB might be more than the available CPU:12 and memory :6.331409454345703GB. This Op [video_split_by_duration_mapper] might require more resource to run.
2024-07-04 14:41:11 | INFO   | data_juicer.core.data:162 - OP [video_split_by_duration_mapper] Done in 18.587(s). Left 12 samples.
2024-07-04 14:41:11 | INFO   | data_juicer.core.data:167 - All OPs are done in 28.465(s).
2024-07-04 14:41:11 | INFO   | data_juicer.core.executor:168 - Exporting dataset to disk...
2024-07-04 14:41:11 | INFO   | data_juicer.core.exporter:140 - Export dataset into a single file...
Creating json from Arrow format: 100% | 1/1 [00:00<00:00, 332.78ba/s]
```

Advanced Processing Hands-on

- Jupyter notebook part 4
 - how to [efficiently process data](#) with CUDA, OP fusion, ...

Tutorial Outline

- **Foundational Abilities**
 - Building Blocks of Data Processing: Data-Juicer's Operators
 - Composing Atomic Capabilities: Data-Juicer's Data Recipes
- **Advanced Data Processing**
 - Exploring Data Recipes: The Data-Juicer Sandbox Lab
 - From Exploration to Production: High-Performance Data Factory
- **Use Cases: From Text to Video Data Processing**
- **Resources and Conclusion**

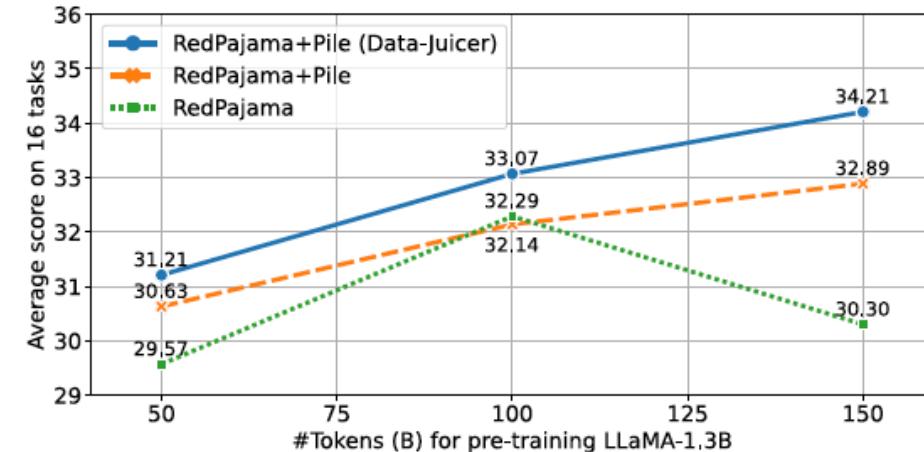
Built-in Text Data Recipe

- How we build the data recipe?
 - Analyze all stats on the whole dataset
 - Set hyper-parameters of OPs with 3-sigma strategy
- Evaluation of the recipe: experiment settings

Exp Item	Pretraining	Finetuning (EN)	Finetuning (CN)
Base Model	LLaMA-1.3B	LLaMA-7B	LLaMA2-7B
Baseline Datasets	RedPajama & The Pile	Alpaca	Belle
Eval Benchmark	16 tasks in HELM	GPT-4 scoring	

Built-in Text Data Recipe

- Pretrain Performance of processed datasets
 - Consistently better performance
 - ✓ Refined datasets > Original datasets
- SFT Performance of processed datasets
 - English model:
 - ✓ Less than 4/5 samples → 17.5% higher winning rate
 - Chinese model:
 - ✓ Less than 1/10 samples → slightly higher winning rate



Model	Tuning Data	#Samples	Win	Tie
LLaMA-7B [33]	Alpaca	52k	16	100
	Data-Juicer	40k	44	
	Random (CFT, EN)	40k	19	105
	Data-Juicer	40k	36	
LLaMA2-7B (Chinese, FlagAlpha [41])	Belle	543k	28	99
	Data-Juicer	52k	33	
	Random (CFT, ZH)	52k	19	96
	Data-Juicer	52k	45	

Built-in Image & Video Data Recipes

- How we build the data recipes?
 - Quick experiments on single/multiple OPs via Data-Juicer Sandbox
 - Apply insights from sandbox:
 - ✓ Select top OPs and combine them, then repeatedly use good data
- Evaluation of the recipes: experiment settings

Exp Item	Image-to-text Task	Text-to-Video Task
Base Model	MGM-2B	T2V-Turbo
Baseline Datasets	MGM pretrain data (LLaVA + ALLaVA-4V)	Panda-70M, InternVid, MSR-VTT
Eval Benchmark	TextVQA, MMBench, MME	V-Bench

Built-in Image Data Recipe

- Performance of processed image dataset
 - Compared to original pre-train dataset: 1/10 distinct instances, 1/2 training instances
 - **Better performance** than baseline model, while **fewer data size (computation)**!

MGM-2B	Num. of Instances	Avg. Perf. Changes (%)	TextVQA	MMBench	MMBench-CN	MME-Perception	MME-Cognition
Official	1226k	-	56.2	59.8	-	1341	312
Reproduced	1226k	-	56.2	59.2	51.6	1334	302
Data-Juicer	159k (x4)	+1.06	54.41	62.08	52.32	1323.37	311.07

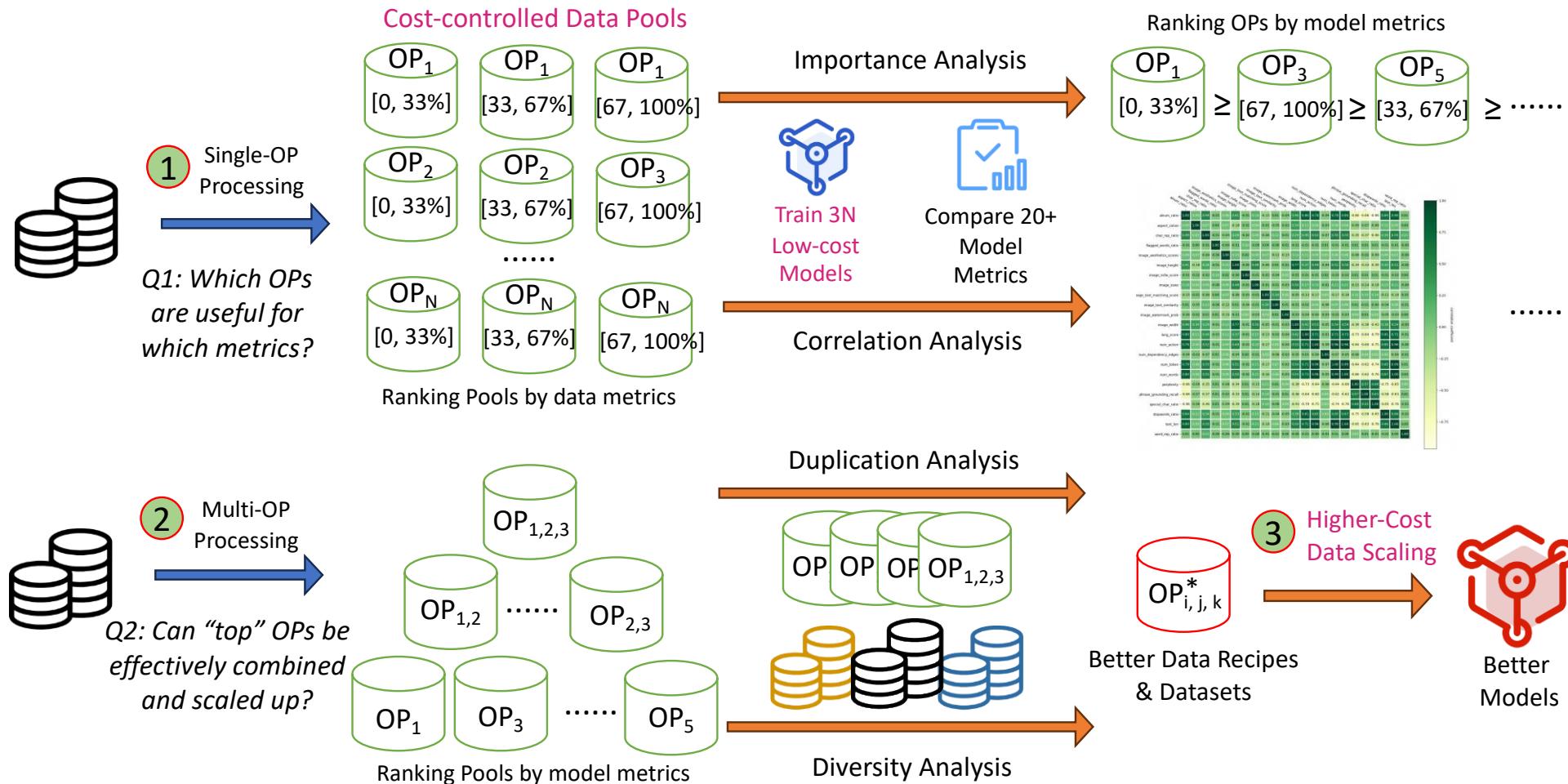
(repeat 4 times)

Built-in Video Data Recipe

- Performance of processed video datasets
 - Thanks to the probe-analyze-refine workflow in Data-Juicer Sandbox
 - Slight processing on data, **achieving a new SOTA performance!**

Models (Ranked by leaderboard)	Board Avg. (%)	Uniform Avg. (%)	Quality Avg. (%)	Semantic Avg. (%)
1. Data-Juicer (T2V)	82.10	80.54	83.14	77.93
2. VEnhancer (VC2) (He et al., 2024a)	81.97	80.00	83.27	76.73
3. LaVie-2 (Wang et al., 2023a)	81.75	79.50	83.24	75.76
4. T2V-Turbo (VC2) (Li et al., 2024)	81.01	78.67	82.57	74.76
5. Gen-2 (Esser et al., 2023)	80.58	80.11	82.47	77.75
6. VideoCrafter-2.0 (Chen et al., 2024b)	80.44	77.81	82.20	73.42
7. Pika Beta (2023-06) (Pika)	80.40	76.97	82.68	71.26
8. AnimateDiff-V2 (Guo et al., 2024)	80.27	76.33	82.90	69.75

Recap: Data-Juicer Sandbox



Sandbox Use Case: Studied OPs

OP Name	Modality	Statistics	Description
alphanumeric_filter	text	Alphanumeric Ratio	Alphanumeric ratio in the text.
character_repetition_filter	text	Character Repetition Ratio	Char-level n-gram repetition ratio in text.
flagged_words_filter	text	Flagged Word Ratio	Flagged-word ratio in the text
image_aesthetics_filter	image	Image Aesthetics Score	Aesthetics score of the image
image_aspect_ratio_filter	image	Image Aspect Ratio	Aspect ratio of the image
image_nsfw_filter	image	Image NSFW Score	NSFW score of the image
image_shape_filter	image	Image Width/Height	Width and height of the image
image_size_filter	image	Image Size	Size in bytes of the image
image_text_matching_filter	text, image	BLIP Image-Text Similarity	Image-text classification matching score based on a BLIP model
image_text_similarity_filter	text, image	CLIP Image-Text Similarity	Image-text feature cosine similarity based on a CLIP model
image_watermark_filter	image	Image Watermark Score	Predicted watermark score of the image based on an image classification model
language_id_score_filter	text	Language Score	Predicted confidence score of the specified language
perplexity_filter	text	Text Perplexity	Perplexity score of the text
phrase_grounding_recall_filter	text, image	Phrase Grounding Recall	Locating recall of phrases extracted from text in the image
special_characters_filter	text	Special Character Ratio	Special character ratio in the text
stopwords_filter	text	Stopword Ratio	Stopword ratio in the text
text_action_filter	text	Text Action Number	Number of actions in the text

Sandbox Use Case: Studied OPs

text_entity_dependency_filter	text	Entity Dependency Number	Number of dependency edges for an entity in the dependency tree of the text
text_length_filter	text	Text Length	Length of the text
token_num_filter	text	Token Number	Token number of the text
video_aesthetics_filter	video	Video Aesthetics Score	Aesthetics score of sampled frames in the video
video_aspect_ratio_filter	video	Video Aspect Ratio	Aspect ratio of the video
video_duration_filter	video	Video Duration	Duration of the video
video_frames_text_similarity_filter	text, video	Frames-Text Similarity	Similarities between sampled frames and text based on a CLIP/BLIP model
video_motion_score_filter	video	Video Motion Score	Motion score of the video
video_nsfw_filter	video	Video NSFW Score	NSFW score of the video
video_ocr_area_ratio_filter	video	Video OCR-Area Ratio	Detected text area ratio for sampled frames in the video
video_resolution_filter	video	Video Width/Height	Width and height of the video
video_watermark_filter	video	Video Watermark Score	Predicted watermark score of the sampled frames in the video based on an image classification model
words_num_filter	text	Word Number	Number of words in the text
word_repetition_filter	text	Word Repetition Ratio	Word-level n-gram repetition ratio in the text

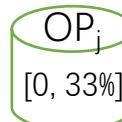
Sandbox: Single-OP case

(compared to random & equal-sized data pool)

Task	OP-Generated Statistics	Average Performance Changes (%)		
		Data Pool (Low)	Data Pool (Mid)	Data Pool (High)
Image-to-Text	Image NSFW Score	7.13 ± 4.29	18.44 ± 18.45	66.38 ± 32.65
	Text Action Number	59.90 ± 46.49	0.29 ± 2.16	-2.05 ± 2.48
	Language Score	49.90 ± 53.82	0.85 ± 2.87	-1.43 ± 2.40
	CLIP Image-Text Similarity	1.20 ± 4.86	-1.81 ± 2.88	49.81 ± 44.72
	Phrase Grounding Recall	-0.49 ± 3.87	-0.58 ± 6.12	49.39 ± 29.83
	Image Width	42.04 ± 57.27	10.31 ± 12.59	1.35 ± 4.36
	Special Character Ratio	-3.08 ± 0.63	-0.75 ± 1.61	39.67 ± 58.82
	Flagged Word Ratio	38.48 ± 27.76	-0.39 ± 0.43	22.49 ± 29.81
	Image Height	35.66 ± 48.62	12.91 ± 10.42	18.73 ± 27.32
	Word Repetition Ratio	33.14 ± 23.39	2.59 ± 5.31	-0.55 ± 2.90
	Aesthetics Score	11.94 ± 12.21	16.58 ± 25.70	0.16 ± 3.67



(ranked by model performance changes)



Sandbox: Single-OP case

Task	OP-Generated Statistics	Average Performance Changes (%)		
		Data Pool (Low)	Data Pool (Mid)	Data Pool (High)
Text-related	Image NSFW Score	7.13 ± 4.29	18.44 ± 18.45	66.38 ± 32.65
	Text Action Number	59.90 ± 46.49	0.29 ± 2.16	-2.05 ± 2.48
	Language Score	49.90 ± 53.82	0.85 ± 2.87	-1.43 ± 2.40
Image-to-Text	CLIP Image-Text Similarity	1.20 ± 4.86	-1.81 ± 2.88	49.81 ± 44.72
	Phrase Grounding Recall	-0.49 ± 3.87	-0.58 ± 6.12	49.39 ± 29.83
	Image Width	42.04 ± 57.27	10.31 ± 12.59	1.35 ± 4.36
	Special Character Ratio	-3.08 ± 0.63	-0.75 ± 1.61	39.67 ± 58.82
	Flagged Word Ratio	38.48 ± 27.76	-0.39 ± 0.43	22.49 ± 29.81
	Image Height	35.66 ± 48.62	12.91 ± 10.42	18.73 ± 27.32
	Word Repetition Ratio	33.14 ± 23.39	2.59 ± 5.31	-0.55 ± 2.90
	Aesthetics Score	11.94 ± 12.21	16.58 ± 25.70	0.16 ± 3.67

Observation 1 (data v.s. model)

Generative models' efficacy is intimately tied to the fidelity of modalities they are trained to generate, which can be explicitly reflected in filtering processes on training data.

Sandbox: Single-OP case

Task	OP-Generated Statistics	Average Performance Changes (%)		
		Data Pool (Low)	Data Pool (Mid)	Data Pool (High)
Video-related	Video Aesthetics Score	-0.98 ± 0.08	0.13 ± 0.09	0.96 ± 0.13
	Video NSFW Score	0.82 ± 0.36	-0.05 ± 0.07	-0.57 ± 0.07
	Frames-Text Similarity	-1.45 ± 0.69	0.23 ± 0.21	0.79 ± 0.15
Text-to-Video	Special-Characters Ratio	0.54 ± 0.36	-0.13 ± 0.70	-0.14 ± 0.10
	Token Number	0.53 ± 0.04	0.18 ± 0.32	0.41 ± 0.25
	Video Height	-0.10 ± 0.21	0.12 ± 0.13	0.46 ± 0.44
	Video OCR-Area Ratio	0.44 ± 0.04	0.02 ± 0.63	-0.66 ± 0.23
	Word Number	-0.49 ± 0.07	-0.41 ± 0.72	0.44 ± 0.45
	Text Action Number	0.18 ± 0.56	-0.71 ± 0.28	0.37 ± 0.28
	Video Motion Score	-0.55 ± 0.40	0.33 ± 0.21	0.32 ± 0.15
	Video Aspect Ratio	-0.32 ± 0.14	0.11 ± 0.18	-0.02 ± 0.40
	Language Score	-0.21 ± 0.01	-0.03 ± 0.38	0.09 ± 0.03
	Video Duration	-0.58 ± 0.05	-0.16 ± 0.09	0.04 ± 0.84

Observation 1 (data v.s. model)

Generative models' efficacy is intimately tied to the fidelity of modalities they are trained to generate, which can be explicitly reflected in filtering processes on training data.

Sandbox: Single-OP case

Task	OP-Generated Statistics	Average Performance Changes (%)		
		Data Pool (Low)	Data Pool (Mid)	Data Pool (High)
Image-to-Text	<u>Image NSFW Score</u>	7.13 ± 4.29	18.44 ± 18.45	66.38 ± 32.65
	Text Action Number	59.90 ± 46.49	0.29 ± 2.16	-2.05 ± 2.48
	<u>Language Score</u>	49.90 ± 53.82	0.85 ± 2.87	-1.43 ± 2.40
	CLIP Image-Text Similarity	1.20 ± 4.86	-1.81 ± 2.88	49.81 ± 44.72
	Phrase Grounding Recall	-0.49 ± 3.87	-0.58 ± 6.12	49.39 ± 29.83
	Image Width	42.04 ± 57.27	10.31 ± 12.59	1.35 ± 4.36
	Special Character Ratio	-3.08 ± 0.63	-0.75 ± 1.61	39.67 ± 58.82
	Flagged Word Ratio	38.48 ± 27.76	-0.39 ± 0.43	22.49 ± 29.81
	Image Height	35.66 ± 48.62	12.91 ± 10.42	18.73 ± 27.32
	Word Repetition Ratio	33.14 ± 23.39	2.59 ± 5.31	-0.55 ± 2.90
	<u>Aesthetics Score</u>	11.94 ± 12.21	16.58 ± 25.70	0.16 ± 3.67

Observation 2

(diversity v.s. quality)

Image-to-text models place greater emphasis on data diversity, whereas text-to-video models prioritize data quality.

Sandbox: Single-OP case

Task	OP-Generated Statistics	Average Performance Changes (%)		
		Data Pool (Low)	Data Pool (Mid)	Data Pool (High)
Text-to-Video	<u>Video Aesthetics Score</u>	-0.98 ± 0.08	0.13 ± 0.09	0.96 ± 0.13
	<u>Video NSFW Score</u>	0.82 ± 0.36	-0.05 ± 0.07	-0.57 ± 0.07
	Frames-Text Similarity	-1.45 ± 0.69	0.23 ± 0.21	0.79 ± 0.15
	Special-Characters Ratio	0.54 ± 0.36	-0.13 ± 0.70	-0.14 ± 0.10
	Token Number	0.53 ± 0.04	0.18 ± 0.32	0.41 ± 0.25
	Video Height	-0.10 ± 0.21	0.12 ± 0.13	0.46 ± 0.44
	Video OCR-Area Ratio	0.44 ± 0.04	0.02 ± 0.63	-0.66 ± 0.23
	Word Number	-0.49 ± 0.07	-0.41 ± 0.72	0.44 ± 0.45
	Text Action Number	0.18 ± 0.56	-0.71 ± 0.28	0.37 ± 0.28
	Video Motion Score	-0.55 ± 0.40	0.33 ± 0.21	0.32 ± 0.15
	Video Aspect Ratio	-0.32 ± 0.14	0.11 ± 0.18	-0.02 ± 0.40
	<u>Language Score</u>	-0.21 ± 0.01	-0.03 ± 0.38	0.09 ± 0.03
	Video Duration	-0.58 ± 0.05	-0.16 ± 0.09	0.04 ± 0.84

Observation 2

(diversity v.s. quality)

Image-to-text models place greater emphasis on data diversity, whereas text-to-video models prioritize data quality.

Sandbox: Single-OP case

Task	OP-Generated Statistics	Average Performance Changes (%)		
		Data Pool (Low)	Data Pool (Mid)	Data Pool (High)
Image-to-Text	Image NSFW Score	7.13 ± 4.29	18.44 ± 18.45	66.38 ± 32.65
	Text Action Number	59.90 ± 46.49	0.29 ± 2.16	-2.05 ± 2.48
	Language Score	49.90 ± 53.82	0.85 ± 2.87	-1.43 ± 2.40
	CLIP Image-Text Similarity	1.20 ± 4.86	-1.81 ± 2.88	49.81 ± 44.72
	Phrase Grounding Recall	-0.49 ± 3.87	-0.58 ± 6.12	49.39 ± 29.83
	Image Width	42.04 ± 57.27	10.31 ± 12.59	1.35 ± 4.36
	Special Character Ratio	-3.08 ± 0.63	-0.75 ± 1.61	39.67 ± 58.82
	Flagged Word Ratio	38.48 ± 27.76	-0.39 ± 0.43	22.49 ± 29.81
	Image Height	35.66 ± 48.62	12.91 ± 10.42	18.73 ± 27.32
	Word Repetition Ratio	33.14 ± 23.39	2.59 ± 5.31	-0.55 ± 2.90
	Aesthetics Score	11.94 ± 12.21	16.58 ± 25.70	0.16 ± 3.67

Observation 3

(spatiotemporal dynamics)

Dynamic information in the data poses a heightened learning challenge for image-to-text generation compared to text-to-video generation.

Sandbox: Single-OP case

Task	OP-Generated Statistics	Average Performance Changes (%)		
		Data Pool (Low)	Data Pool (Mid)	Data Pool (High)
Text-to-Video	Video Aesthetics Score	-0.98 ± 0.08	0.13 ± 0.09	0.96 ± 0.13
	Video NSFW Score	0.82 ± 0.36	-0.05 ± 0.07	-0.57 ± 0.07
	Frames-Text Similarity	-1.45 ± 0.69	0.23 ± 0.21	0.79 ± 0.15
	Special-Characters Ratio	0.54 ± 0.36	-0.13 ± 0.70	-0.14 ± 0.10
	Token Number	0.53 ± 0.04	0.18 ± 0.32	0.41 ± 0.25
	Video Height	-0.10 ± 0.21	0.12 ± 0.13	0.46 ± 0.44
	Video OCR-Area Ratio	0.44 ± 0.04	0.02 ± 0.63	-0.66 ± 0.23
	Word Number	-0.49 ± 0.07	-0.41 ± 0.72	0.44 ± 0.45
	Text Action Number	0.18 ± 0.56	-0.71 ± 0.28	0.37 ± 0.28
	Video Motion Score	-0.55 ± 0.40	0.33 ± 0.21	0.32 ± 0.15
Image-to-Text	Video Aspect Ratio	-0.32 ± 0.14	0.11 ± 0.18	-0.02 ± 0.40
	Language Score	-0.21 ± 0.01	-0.03 ± 0.38	0.09 ± 0.03
	Video Duration	-0.58 ± 0.05	-0.16 ± 0.09	0.04 ± 0.84

Observation 3

(spatiotemporal dynamics)

Dynamic information in the data poses a heightened learning challenge for image-to-text generation compared to text-to-video generation.

Sandbox: Single-OP case

Task	OP-Generated Statistics	Average Performance Changes (%)		
		Data Pool (Low)	Data Pool (Mid)	Data Pool (High)
Image-to-Text	Image NSFW Score	7.13 ± 4.29	18.44 ± 18.45	66.38 ± 32.65
	Text Action Number	59.90 ± 46.49	0.29 ± 2.16	-2.05 ± 2.48
	Language Score	49.90 ± 53.82	0.85 ± 2.87	-1.43 ± 2.40
	CLIP Image-Text Similarity	1.20 ± 4.86	-1.81 ± 2.88	49.81 ± 44.72
	Phrase Grounding Recall	-0.49 ± 3.87	-0.58 ± 6.12	49.39 ± 29.83
	Image Width	42.04 ± 57.27	10.31 ± 12.59	1.35 ± 4.36
	Special Character Ratio	-3.08 ± 0.63	-0.75 ± 1.61	39.67 ± 58.82
	Flagged Word Ratio	38.48 ± 27.76	-0.39 ± 0.43	22.49 ± 29.81
	Image Height	35.66 ± 48.62	12.91 ± 10.42	18.73 ± 27.32
	Word Repetition Ratio	33.14 ± 23.39	2.59 ± 5.31	-0.55 ± 2.90
	Aesthetics Score	11.94 ± 12.21	16.58 ± 25.70	0.16 ± 3.67

Observation 4

(modality alignment)

A high degree of match between different modalities within the data is crucial for model performance in both image-to-text and text-to-video generation.

Sandbox: Single-OP case

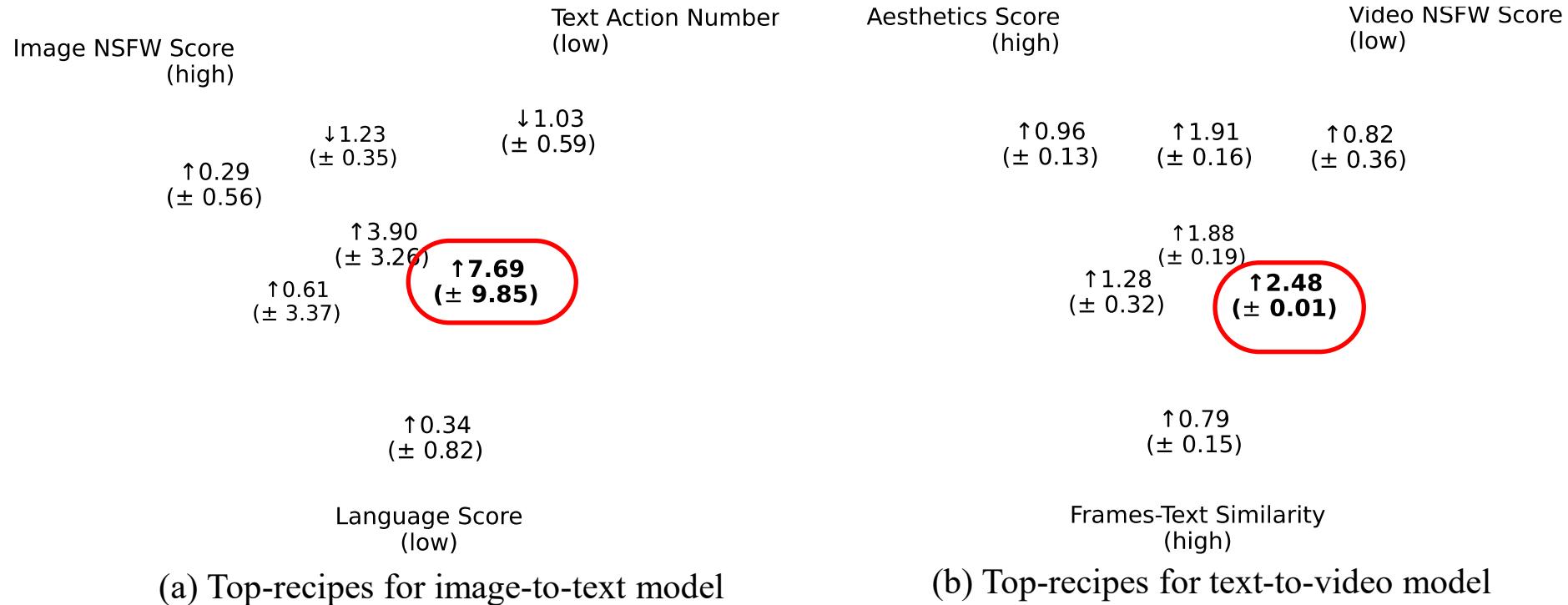
Task	OP-Generated Statistics	Average Performance Changes (%)		
		Data Pool (Low)	Data Pool (Mid)	Data Pool (High)
Text-to-Video	Video Aesthetics Score	-0.98 ± 0.08	0.13 ± 0.09	0.96 ± 0.13
	Video NSFW Score	0.82 ± 0.36	-0.05 ± 0.07	-0.57 ± 0.07
	Frames-Text Similarity	-1.45 ± 0.69	0.23 ± 0.21	0.79 ± 0.15
	Special-Characters Ratio	0.54 ± 0.36	-0.13 ± 0.70	-0.14 ± 0.10
	Token Number	0.53 ± 0.04	0.18 ± 0.32	0.41 ± 0.25
	Video Height	-0.10 ± 0.21	0.12 ± 0.13	0.46 ± 0.44
	Video OCR-Area Ratio	0.44 ± 0.04	0.02 ± 0.63	-0.66 ± 0.23
	Word Number	-0.49 ± 0.07	-0.41 ± 0.72	0.44 ± 0.45
	Text Action Number	0.18 ± 0.56	-0.71 ± 0.28	0.37 ± 0.28
	Video Motion Score	-0.55 ± 0.40	0.33 ± 0.21	0.32 ± 0.15
	Video Aspect Ratio	-0.32 ± 0.14	0.11 ± 0.18	-0.02 ± 0.40
	Language Score	-0.21 ± 0.01	-0.03 ± 0.38	0.09 ± 0.03
	Video Duration	-0.58 ± 0.05	-0.16 ± 0.09	0.04 ± 0.84

Observation 4

(modality alignment)

A high degree of match between different modalities within the data is crucial for model performance in both image-to-text and text-to-video generation.

Sandbox: Multi-OPs case



Observation 5

($1 + 1 < 2$, sequentially)

The optimal data recipe does not necessarily result from combining the best individual OPs, nor does incorporating more high-performing OPs always lead to superior outcomes.

Sandbox: Multi-OPs case

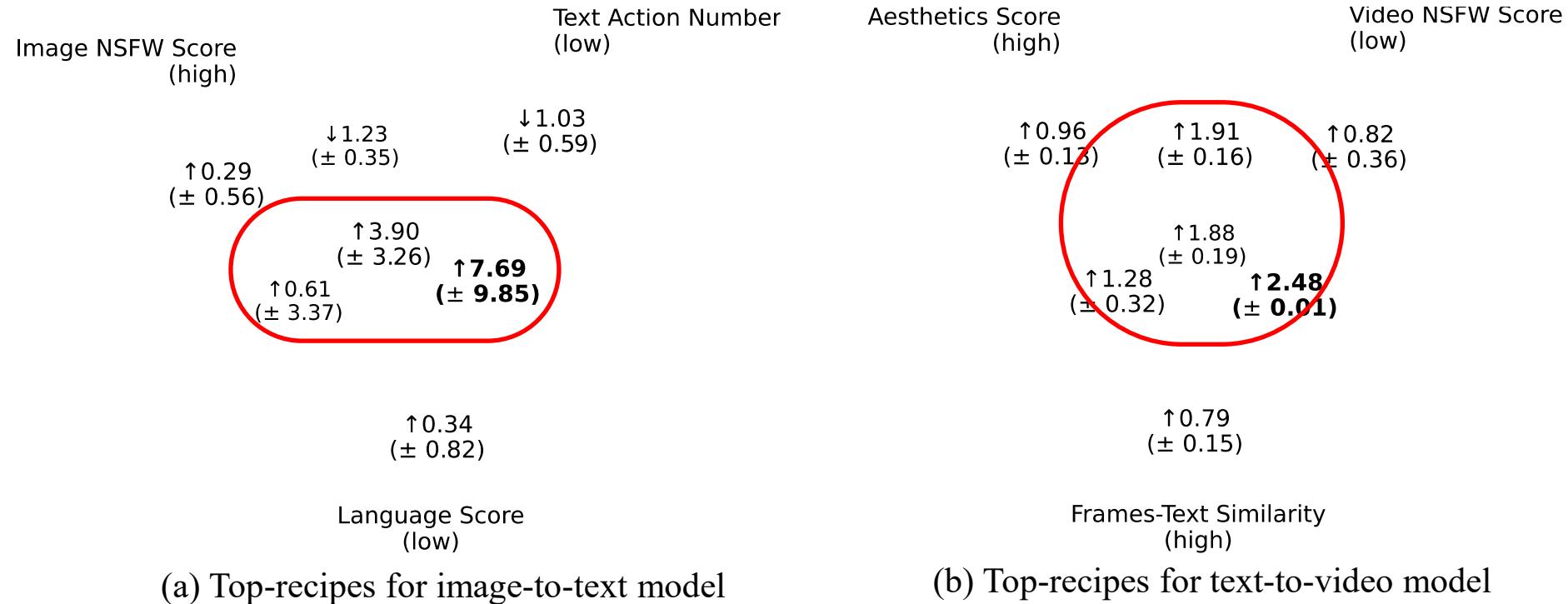
Text Action Number (low)	Phrase Grounding Recall (high)	Video Motion Score (low)	Aesthetics Score (high)
$\uparrow 5.21$ (± 0.43)	$\uparrow 4.74$ (± 1.35)	$\uparrow 6.09$ (± 2.34)	$\downarrow 0.55$ (± 0.40)
$\uparrow 2.72$ (± 3.84)	$\uparrow 5.54$ (± 2.94)	$\downarrow 1.14$ (± 0.04)	$\uparrow 0.06$ (± 0.76)
$\uparrow 5.40$ (± 2.83)	$\uparrow 3.65$ (± 2.55)	$\downarrow 0.43$ (± 0.77)	$\uparrow 0.33$ (± 0.20)
Image NSFW Score (high)			$\uparrow 0.04$ (± 0.84)
(c) Cluster-recipes for image-to-text model	(d) Cluster-recipes for text-to-video model		

Observation 6

$(1 + 1 < 2, \text{orthogonally})$

Combining OPs that excel in orthogonal dimensions on model or data does not guarantee complementary effects; rather, it is more likely that they will impede each other's performance.

Sandbox: Multi-OPs case



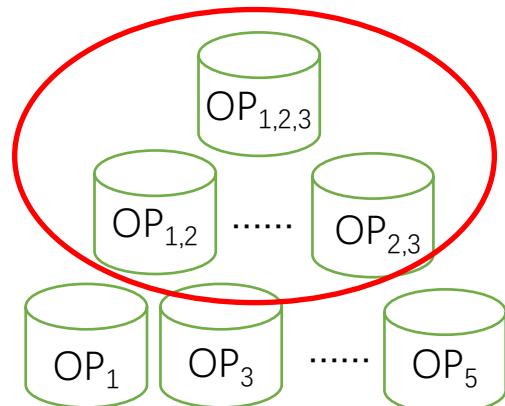
Observation 7

$$(1 + x > x)$$

The performance of a single OP is positively correlated with the performance of the recipe created from its combination. Starting with high-performing OPs is a good initial step in exploring optimal higher-order data recipes.

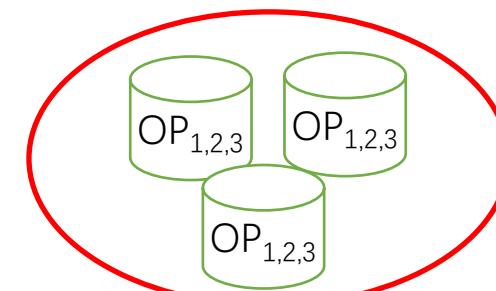
Sandbox: Data Duplication Case

A: Without Duplicates



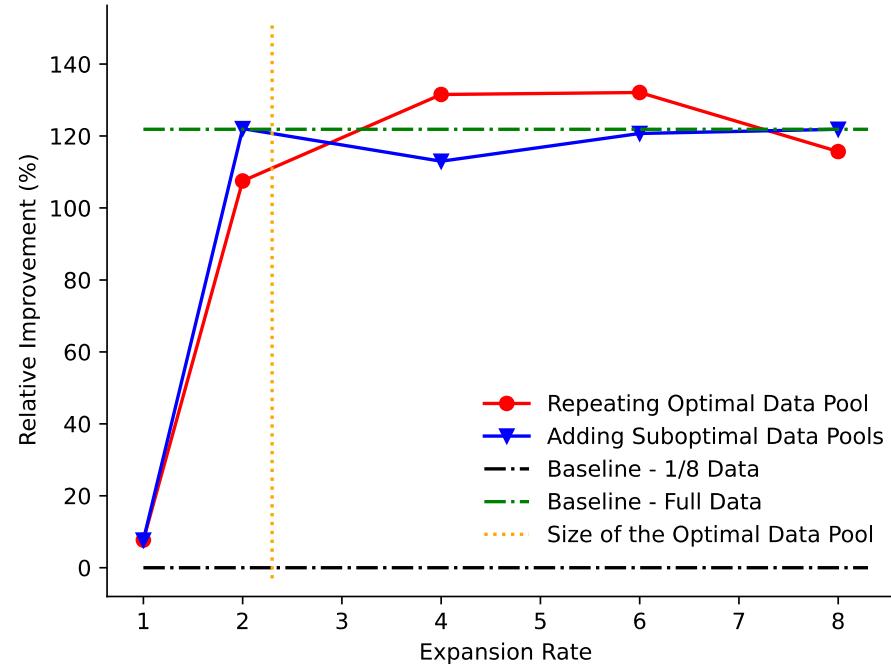
Ranking Pools by model metrics

B: With Duplication

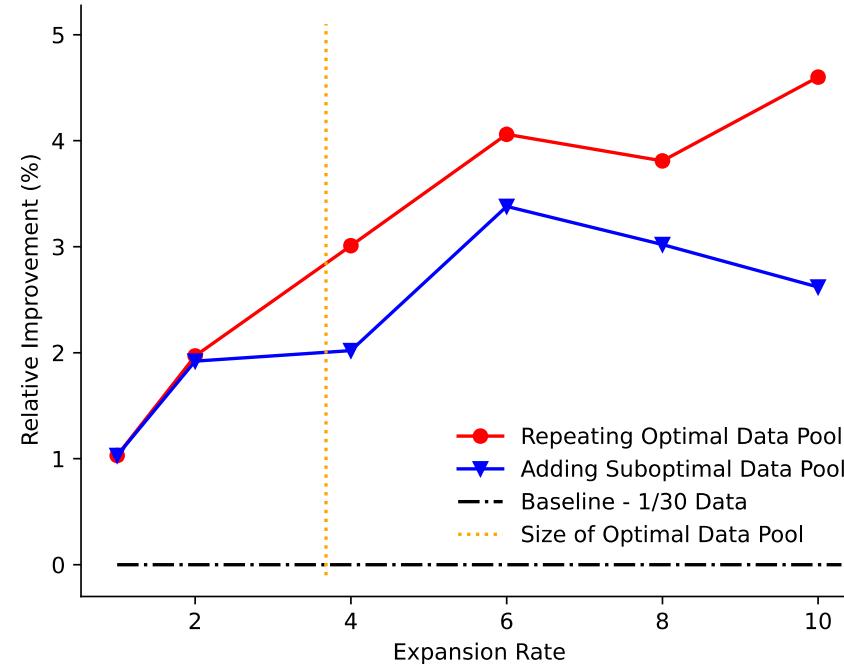


For text-only LLMs, A usually > B, but how about MLLMs?

Sandbox: Data Duplication Case



(a) Image-to-text



(b) Text-to-video

Observation 8

(effect of duplicates)

Duplicating high-quality data is beneficial for both image-to-text and text-to-video models. For the image-to-text model, optimal utilization of high-quality data may be achieved after four repetitions, while for the text-to-video model, it is effective to reuse high-quality data extensively between six to ten times.

The Use-cases Hands-on

- Jupyter notebooks part 5
 - [5.1: Synthesize data](#) for Image-Text datasets
 - [5.2: Process video datasets](#) from spacetime perspectives

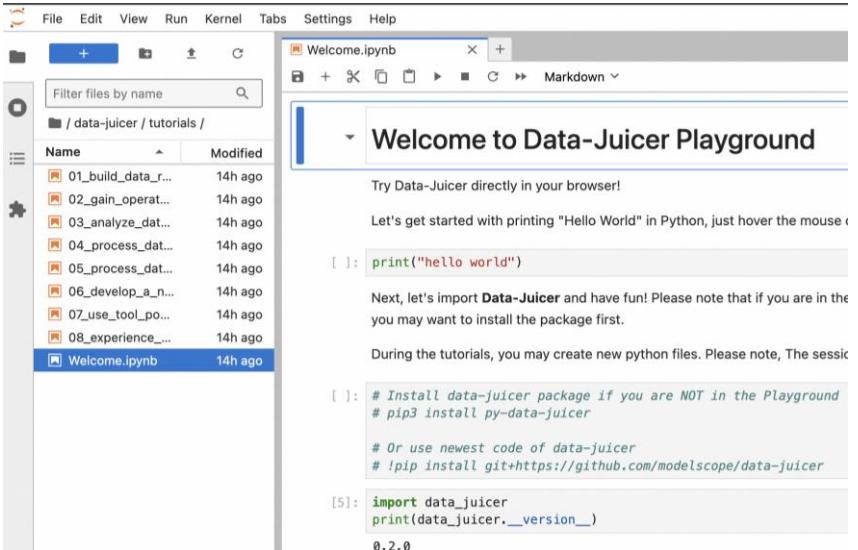
Tutorial Outline

- **Foundational Abilities**
 - Building Blocks of Data Processing: Data-Juicer's Operators
 - Composing Atomic Capabilities: Data-Juicer's Data Recipes
- **Advanced Data Processing**
 - Exploring Data Recipes: The Data-Juicer Sandbox Lab
 - From Exploration to Production: High-Performance Data Factory
- Use Cases: From Text to Video Data Processing
- **Resources and Conclusion**

Interactive Resources

➤ Playgrounds

- Jupyter Notebooks
- Copilot Assistant



The screenshot shows a Jupyter Notebook interface. On the left is a file tree with a folder named 'data-juicer / tutorials /' containing files like '01_build_data_r...', '02_gain_operat...', etc., all modified 14 hours ago. The main area displays a notebook titled 'Welcome.ipynb' with the following content:

```
[ ]: print("hello world")
```

Next, let's import Data-Juicer and have fun! Please note that if you are in the browser, you may want to install the package first.

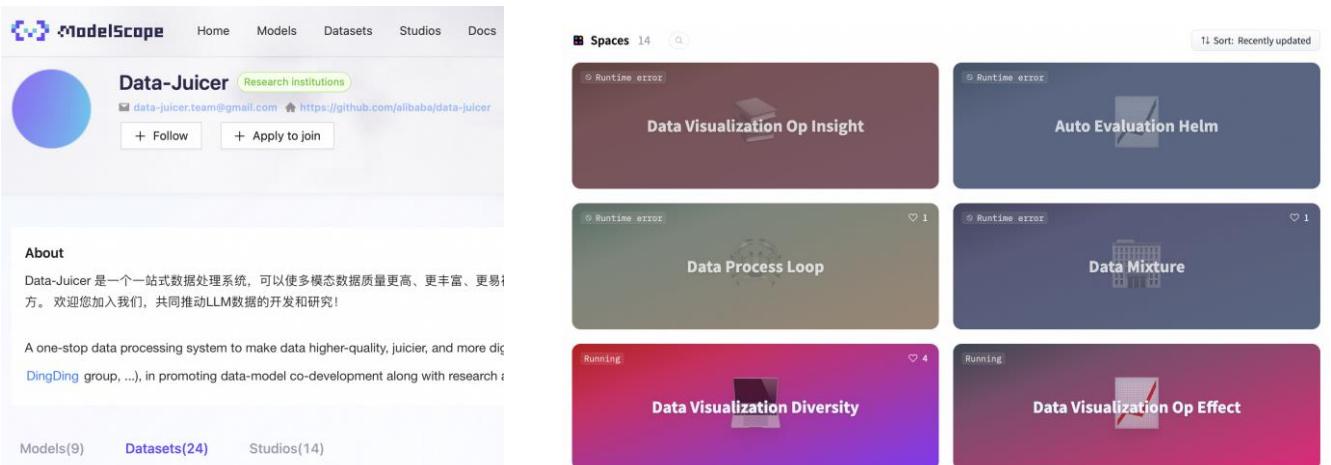
During the tutorials, you may create new python files. Please note, The session will be lost if you leave the browser.

```
[ ]: # Install data-juicer package if you are NOT in the Playground  
# pip3 install py-data-juicer  
  
# Or use newest code of data-juicer  
# !pip install git+https://github.com/modelscope/data-juicer
```

```
[5]: import data_juicer  
print(data_juicer.__version__)
```

0.2.0

- ## ➤ 50+ interactive spaces, models, datasets on ModelScope & HuggingFace

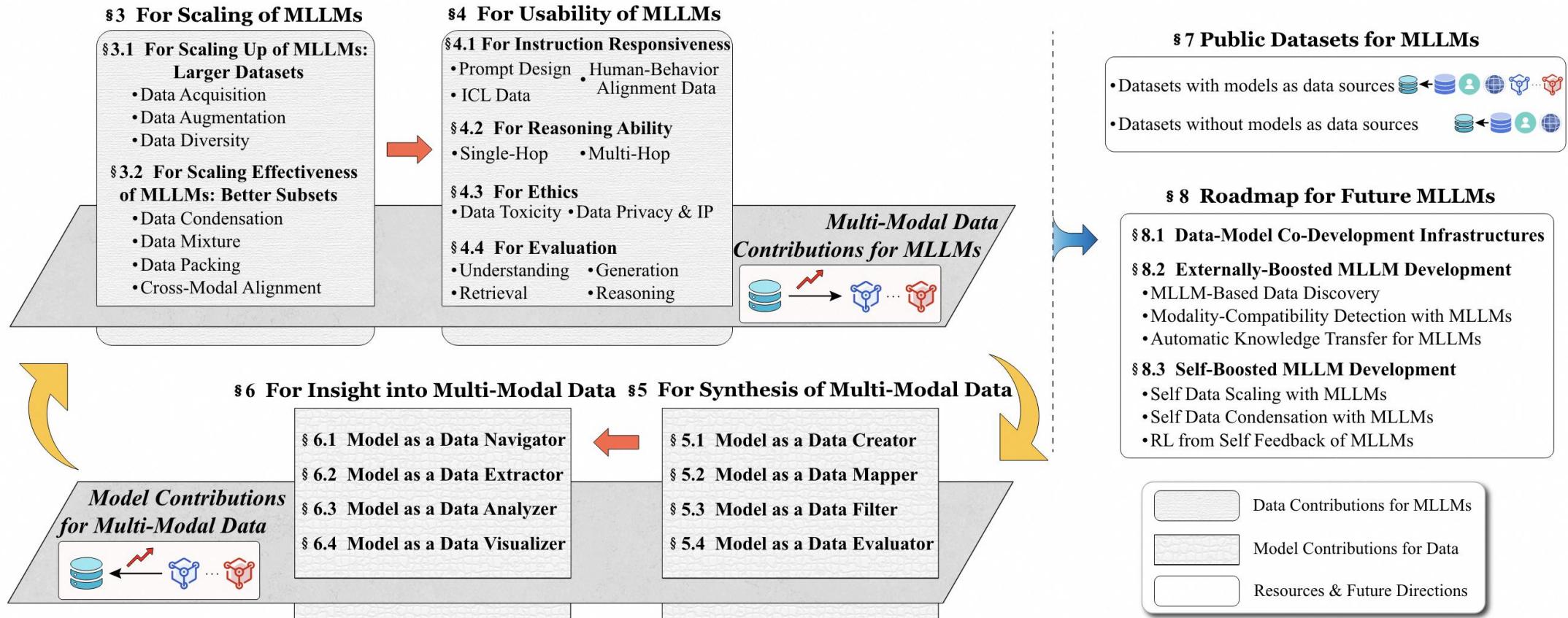


Data-Juicer Competitions

- 2000+ teams, open-sourced solutions
- Historical:
 - FT-Data Ranker: LLM finetuning data cleaning
 - Better Mixture: mix data from various sources
- On-going:
 - ModelScope-SORA: video generation
 - Better Synth: image understanding



The Maintained Survey

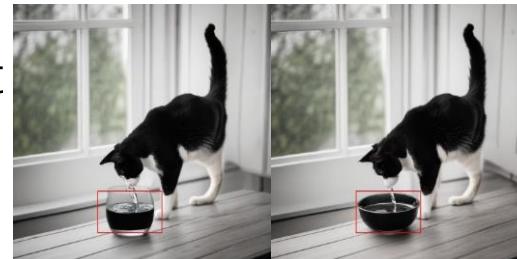


[1] (arXiv:2407.08583) The Synergy between Data and Multi-Modal Large Language Models: A Survey from Co-Development Perspective
https://github.com/modelscope/data-juicer/docs/awesome_llm_data

Future Direction

MLLMs as Data Filter/Mapper

- Alleviate model collapse
 - MLLMs are generative models!
 - Inductive bias propagation → longer tailed distribution
- Better modality-compatibility
 - text-centric → any-to-any alignment
 - contrastive synthesis^[7]
- Automatic compliance
 - privacy
 - licensing requirements



The left image shows a cat drinking water from a **glass**, while the right image shows the same cat drinking water from a **black bowl**. The difference is the type of container used for the water.

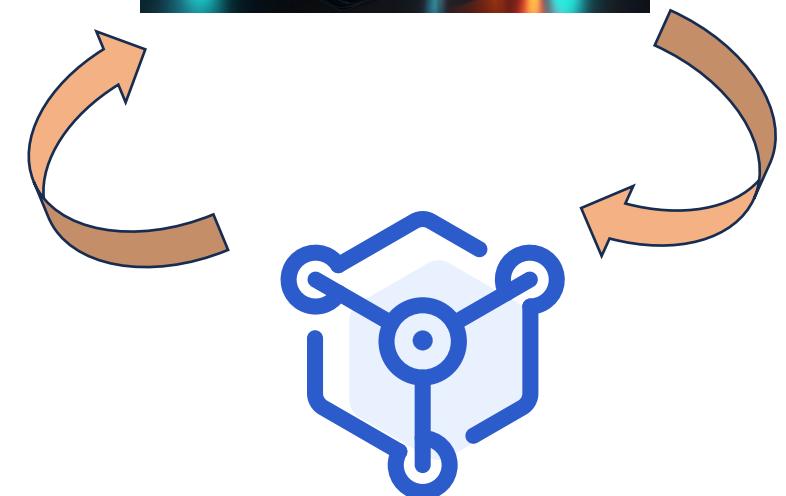


The difference between the two images lies in the object that the player is holding. In the left image, the player is holding a **baseball glove**, while in the right image, the player is holding a **basketball**.

[7] (arXiv:2408.04594) ImgDiff: Contrastive Data Synthesis for Vision Large Language Models

Future Direction Data-Model Self-Boosting

- Dynamic Data & Model
 - Mutual enhancement
 - e.g., [caption₀, fig₀] → train better MLLMs
 - generate/clean better data:
(more detailed, precise [caption₁ , fig₁])
 - train better MLLMs →
- RL from self-feedback of data and MLLMs



Conclusion

➤ Multi-modal Foundation Models

- Scaling Laws (both computation and **data**)
- “Garbage in, garbage out”



➤ Data-Juicer

- open-sourced, effect-proven, industrial system
- data analysis, cleaning, synthesis
- sandbox: data-model co-development

Contributors 24



[+ 10 contributors](#)

➤ Under continuous maintenance, enhancement; welcome to join the community!

References

[1] (arXiv:2407.08583)

The Synergy between Data and Multi-Modal Large Language Models: A Survey from Co-Development Perspective

[2] (SIGMOD'24)

Data-juicer: A one-stop data processing system for large language models

[3] (arXiv:2405.14908)

Data Mixing Made Efficient: A Bivariate Scaling Law for Language Model Pretraining

[4] (arXiv:2407.11784)

Data-Juicer Sandbox: A Comprehensive Suite for Multimodal Data-Model Co-development

[5] (arXiv:2408.04594)

ImgDiff: Contrastive Data Synthesis for Vision Large Language Models

Thanks!



GitHub: <https://github.com/modelscope/data-juicer>