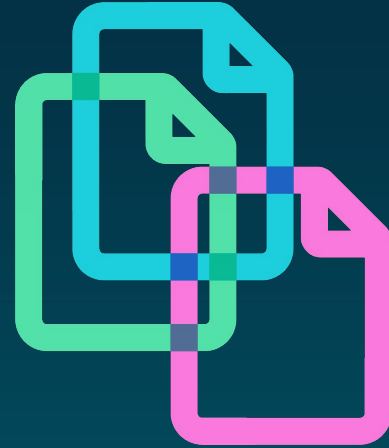# BigDocs

# An Open Dataset for Training Multimodal Models on Document and Code Tasks

*bigdocs.github.io*

**Juan A. Rodriguez, Xiangru Jian, et. al.**
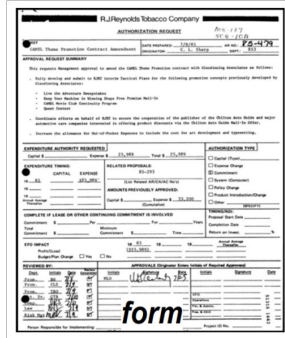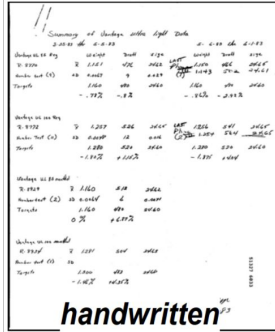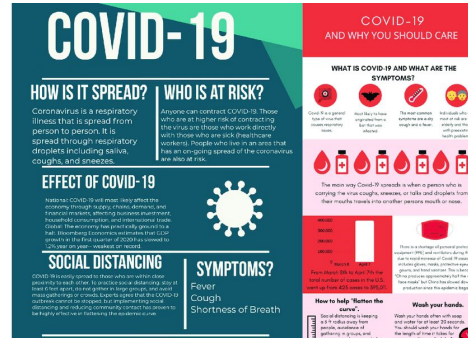**ServiceNow Research**

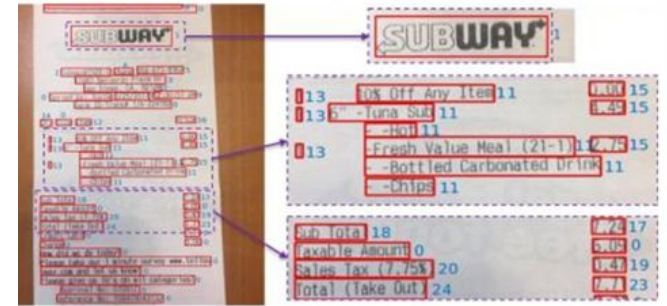juan.rodriguez@mila.quebec
xiangru.jian@waterloo.ca

# Multimodal Document Understanding

# Why Multimodal Document Understanding?

Existing demand of models to deal with:



## Complex Paper Documents



## Text-intensive Infographics
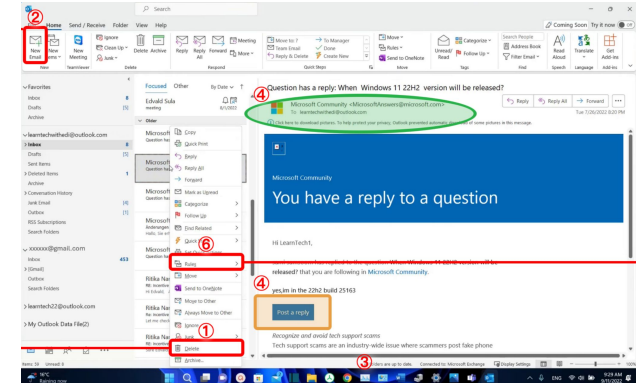


## Robust Perception & Grounding

# Why Multimodal Document Understanding?

New tasks for Document AI



## Multimodal Code Generation



## User Interface Understanding

# Current Limitations:

1. **Models** need *huge* and *diverse* datasets to generalize well.

2. **Data is limited**, scattered, poorly licensed, or simplistic.

3. **Modern/real world tasks** are underexplored

"

*We need **Scalable**, **Open**, and **High-Quality** Multimodal Document Datasets*

# BigDocs Contributions Summary

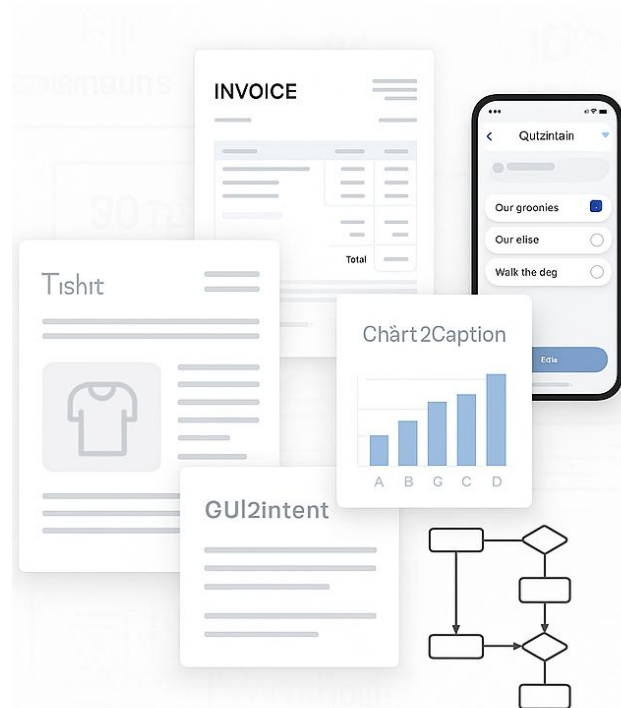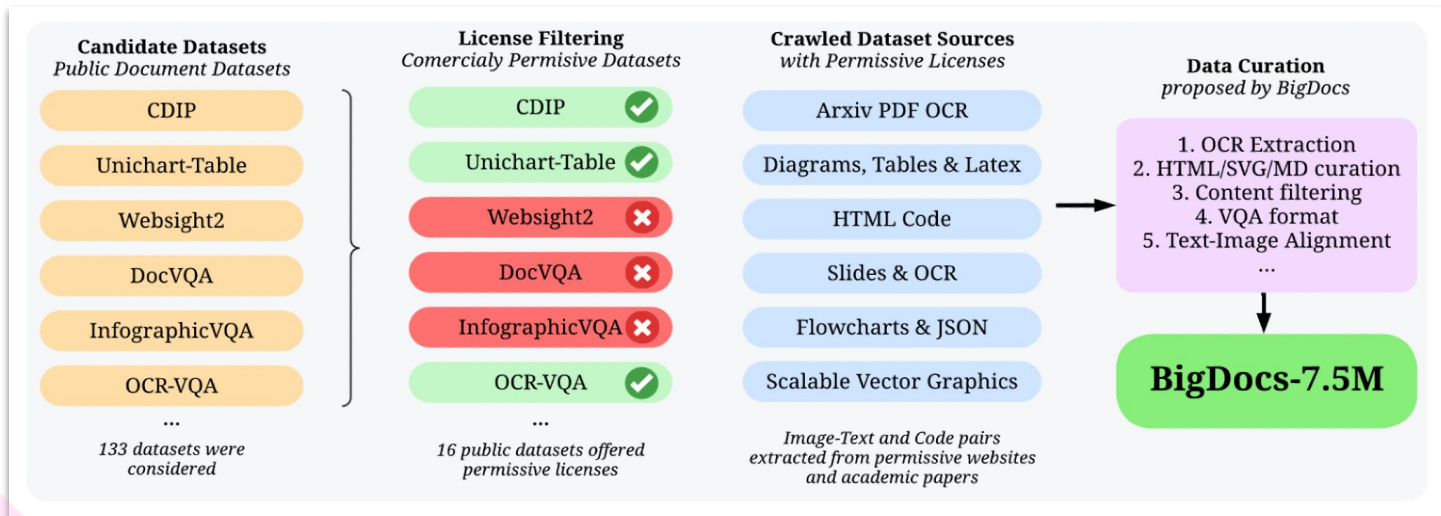| | | |
|---|---|---|
| 🗄️ | **BigDocs 7.5M** | Large-scale pre-training dataset with clearly licensing and transparent |
| 📊 | **BigDocs-Bench** | 10 innovative tasks for multimodal code & GUI Understanding |
| ⚙️ | **BigDocs-Toolkit** | Crawling, dataset compound management, safety filtering, document AI training |
| 🚀 | **BigDocs-Models** | License-permissive multimodal document models  Phi-3  QwenVL2  LLaVA |

servicenow

# Dataset Construction

- Merged and **filtered 130+ datasets** → **kept only 16** with **permissive licenses**

- Added **new data**: UIs, HTML screenshots, LaTeX tables, …

- **Strict content**, **license**, and **quality checks** throughout



**Candidate Datasets**
*Public Document Datasets*

CDIP

Unichart-Table

Websight2

DocVQA

InfographicVQA

OCR-VQA

…
*133 datasets were considered*

**License Filtering**
*Comercialy Permisive Datasets*

CDIP ✅

Unichart-Table ✅

Websight2 ❌

DocVQA ❌

InfographicVQA ❌

OCR-VQA ✅

…
*16 public datasets offered permissive licenses*

**Crawled Dataset Sources**
*with Permissive Licenses*

Arxiv PDF OCR

Diagrams, Tables & Latex

HTML Code

Slides & OCR

Flowcharts & JSON

Scalable Vector Graphics

*Image-Text and Code pairs extracted from permissive websites and academic papers*

**Data Curation**
*proposed by BigDocs*

1. OCR Extraction
2. HTML/SVG/MD curation
3. Content filtering
4. VQA format
5. Text-Image Alignment
…

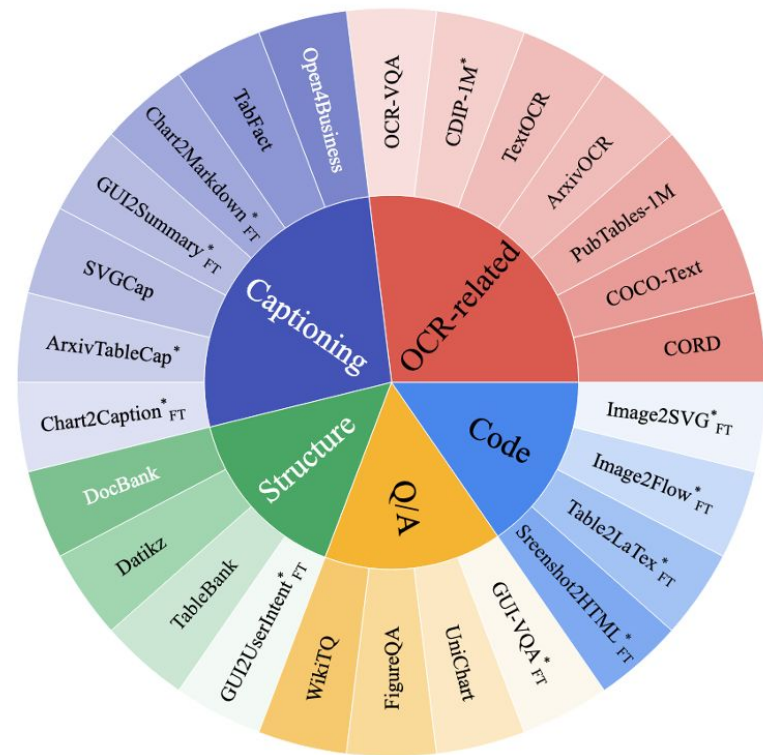**BigDocs-7.5M**

**BigDocs**

# What's in BigDocs-7.5M?

OCR, Layout, Tables, Forms, Charts, UIs, Scientific Diagrams

Text, Images, Structured Outputs (HTML, LaTeX, JSON, SVG)

Balanced across input-output modalities: vision → code, vision → text

# What's in BigDocs-7.5M?

We **curated multimodal datasets** in an **open** and **transparent** way specifically **for training multimodal models document, website, coding domains**

" 

*We created **new Benchmarks** for* **Modern Document AI**

# Introducing **BigDocs-Bench**

- **Proposing 10 new tasks** for challenging real-world use cases
- Tests **visual *reasoning*, layout understanding**, *structured code generation*, *UI comprehension*

**Multimodal Code Generation**

Image2HTML, Image2LaTeX, Image2SVG, Flow2Code

**User Interface Reasoning**

GUI2Intent, GUI2Summary, GUI-VQA

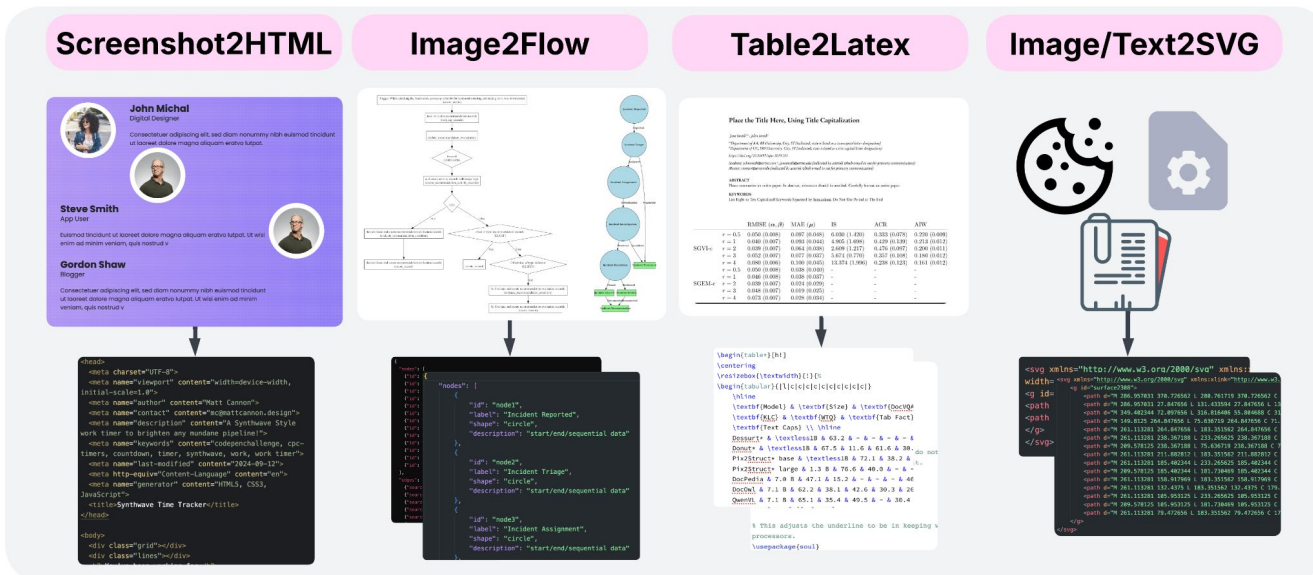**Chart Understanding**

Chart2Caption

"

**It's not just OCR anymore.**

*Models must generate complex outputs from visuals.*

# Novel Tasks in BigDocs-Bench

| Task | Train | Val | Test | Hidden | Tokens |
|------|-------|-----|------|--------|--------|
| 🖥️ Screenshot-2HTML | 9.3K | 1000 | 500 | 500 | 32.7K±53K |
| ⊞ Table-2LaTeX | 77.7K | 1000 | 500 | 500 | 438±540 |
| ⬚ Image2SVG | 198K | 2000 | 748 | 500 | 2.9K±1.7K |
| 🖧 Image2Flow (GraphViz) | 8.0K | 1000 | 500 | 500 | 418±124 |
| </> Image2Flow (JSON) | 8000 | 1000 | 500 | 500 | 1800±601 |
| 📊 Chart-2Markdown | 4500 | 1000 | 500 | 500 | 1.6K±4.4K |
| 📈 Chart2Caption | 5.4K | 1300 | 650 | 500 | 94±49 |
| 👤 GUI2UserIntent | 79K | 1000 | 500 | 500 | 28±4 |
| 📄 GUI2Summary | 79K | 1000 | 500 | 500 | 132±25 |
| ❓ GUI-VQA | 78.9k | 1000 | 500 | 500 | 35±24 |

# Novel Tasks: Image2Code

- Input: Screenshot of web page, table, workflow, logotype
- Output: HTML / LaTeX / SVG
- Requires: understanding *layout*, *content*, and **generating *structured code***
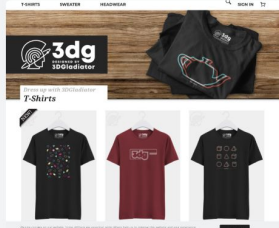
# Novel Tasks: GUI & Chart Understanding

- Answer questions & describe UI screenshots.
- Predict user's action from a UI screenshot
- Interpret data visualizations

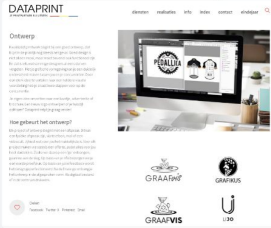# Why use BigDocs?

# Training & Models

- We **Continually Pre-Trained** open-source VLMs to perform **Document AI**

- We **instruction-tuned** the models for multiple downstream tasks

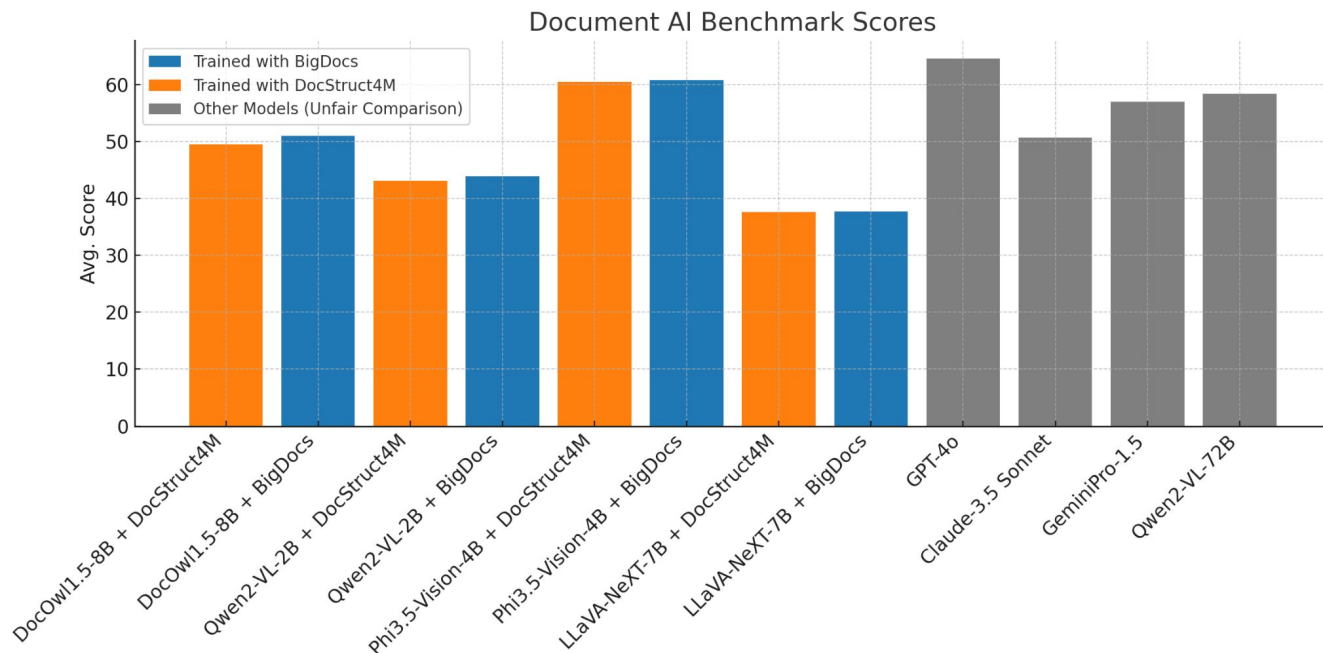- We released full **training/inference** pipelines and recipes

# Results on Document AI Benchmarks

- **Pretraining on BigDocs** surpasses other pre-training approaches
  - No benchmark contamination, transparent, license clarity
- Our best model **matches/outperforms SOTA closed models**



Document AI Benchmark Scores

Legend:
- Trained with BigDocs
- Trained with DocStruct4M
- Other Models (Unfair Comparison)

Y-axis: Avg. Score

X-axis categories: DocOwl1.5-8B + DocStruct4M, DocOwl1.5-8B + BigDocs, Qwen2-VL-2B + DocStruct4M, Qwen2-VL-2B + BigDocs, Phi3.5-Vision-4B + DocStruct4M, Phi3.5-Vision-4B + BigDocs, LLaVA-NeXT-7B + DocStruct4M, LLaVA-NeXT-7B + BigDocs, GPT-4o, Claude-3.5 Sonnet, GeminiPro-1.5, Qwen2-VL-72B
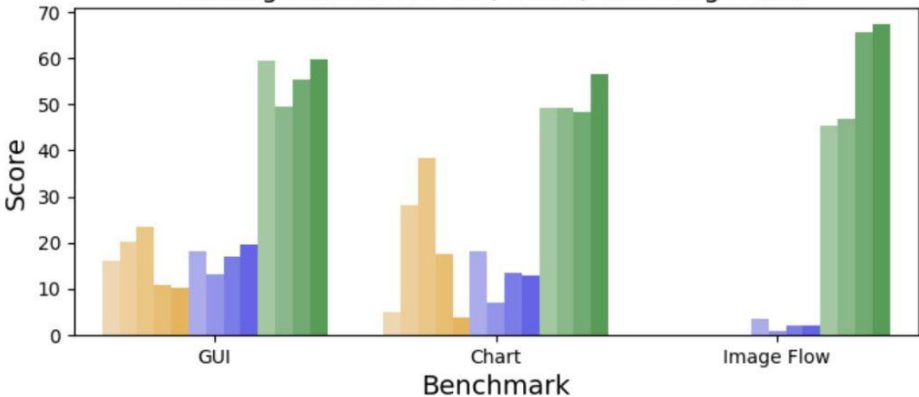
# Results on BigDocs-Bench

Models struggle at Multimodal Coding and GUI Reasoning
BigDocs teaches new skills, surpassing SOTA models



Model Performance Comparison on BigDocs-Bench

# Results - Human Evaluation

- Human evaluation **reveals a clear preference**
- BigDocs **outperforms** GPT4 on **Latex** and **HTML** generation

# Qualitative Example - Latex Generation

- BigDocs achieves a better table conversion than GPT4

**Input Image**

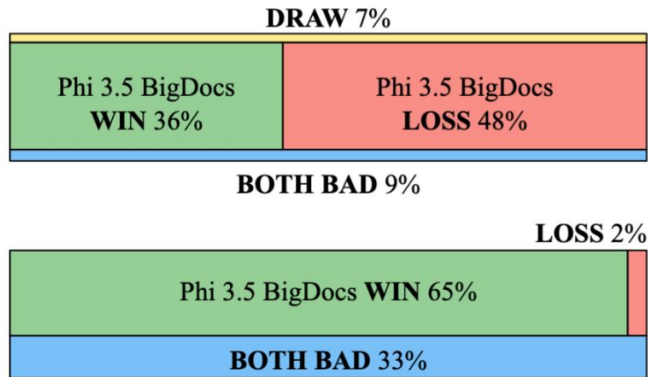|  | Discrete dynamics(ACC) | | | | Continuous dynamics($\sigma$) | | |
|---|---|---|---|---|---|---|---|
|  | SIR | SIS | Threshold | Kirman | Gene | Mutualistic | CML |
| $T+1$ | 0.85 | 0.86 | 0.89 | 0.84 | 0.598 | 0.958 | 0.017 |
| $T+2$ | 0.73 | 0.80 | 0.84 | 0.81 | 0.602 | 1.086 | 0.021 |
| $T+3$ | 0.81 | 0.75 | 0.81 | 0.82 | 0.609 | 1.276 | 0.024 |
| $T+4$ | 0.82 | 0.74 | 0.74 | 0.83 | 0.724 | 1.512 | 0.027 |
| $T+5$ | 0.80 | 0.74 | 0.72 | 0.85 | 0.822 | 1.601 | 0.028 |

**Output from GPT4o**

|  | Discrete dynamics (ACC) | | | | Continuous dynamics ($\sigma$) | | |
|---|---|---|---|---|---|---|---|
|  | SIR | SIS | Threshold | Kirman | Gene | Mutualistic | CML |
| $T+1$ | 0.85 | 0.86 | 0.89 | 0.84 | 0.598 | 0.958 | 0.017 |
| $T+2$ | 0.73 | 0.80 | 0.84 | 0.81 | 0.602 | 1.086 | 0.021 |
| $T+3$ | 0.81 | 0.75 | 0.81 | 0.82 | 0.609 | 1.276 | 0.024 |
| $T+4$ | 0.82 | 0.74 | 0.74 | 0.83 | 0.724 | 1.512 | 0.027 |
| $T+5$ | 0.80 | 0.74 | 0.72 | 0.85 | 0.822 | 1.601 | 0.028 |

**Output from BigDocs**

|  | Discrete dynamics(ACC) | | | | Continuous dynamics($\sigma$) | | |
|---|---|---|---|---|---|---|---|
|  | SIR | SIS | Threshold | Kirman | Gene | Mutualistic | CML |
| $T+1$ | 0.85 | 0.86 | 0.89 | 0.84 | 0.598 | 0.958 | 0.017 |
| $T+2$ | 0.73 | 0.80 | 0.84 | 0.81 | 0.602 | 1.086 | 0.021 |
| $T+3$ | 0.81 | 0.75 | 0.81 | 0.82 | 0.609 | 1.276 | 0.024 |
| $T+4$ | 0.82 | 0.74 | 0.74 | 0.83 | 0.724 | 1.512 | 0.027 |
| $T+5$ | 0.80 | 0.74 | 0.72 | 0.85 | 0.822 | 1.601 | 0.028 |

# Results – Does BigDocs Help?

💪 BigDocs **adapts eneral multimodal models** to reach **SOTA performance in Document AI**.

🏆 Delivers **+34.5% gains** on **code generation** and **GUI understanding** tasks.

🚀 **Outperforms GPT-4** on **LaTeX** generation and intent prediction from **GUIs**.

# Summary & Takeaways

- 📦 **BigDocs-7.5M**: Large open multimodal document dataset

- 🧠 **BigDocs-Bench**: 10 hard tasks, real-world relevance

- 🚀 **Strong gains** – even outperform GPT-4 in key areas

- 🔒 **Fully open** – data, tools, models, benchmarks

# BigDocs

An Open Dataset for Training Multimodal Models on Document and Code Tasks

*bigdocs.github.io*

**Juan A. Rodriguez, Xiangru Jian, et. al.**
**ServiceNow Research**

## Thanks! Contact us

juan.rodriguez@mila.quebec
xiangru.jian@waterloo.ca