# 科技生态圈峰会 + 深度研习

## ——1000 + 技术团队的共同选择

KEYLINK ing

### K+峰会

**上海站**
K+ 峰会
K+ 全球软件研发行业创新峰会
时间：2024.06.21-22

**敦煌站**
K+ 峰会
K+ 思考周®研习社
时间：2024.10.17-19

**香港站**
K+ 峰会
K+ 思考周®研习社
时间：2024.11.10-12

K+峰会详情

### AiDD峰会

**上海站**
AiDD 峰会
AI+研发数字峰会
时间：2024.05.17-18

**北京站**
AiDD 峰会
AI+研发数字峰会
时间：2024.08.16-17

**深圳站**
AiDD 峰会
AI+研发数字峰会
时间：2024.11.08-09

AiDD峰会详情

# AiDD

# 2024 AI+研发数字峰会
## AI+ Development Digital summit
深圳站 11/08-09

# AI驱动研发变革 促进企业降本增效

## 2024深圳站-议题设置

| AI+产品线 | LLM驱动产品创新 | LLM驱动需求与业务分析 | AI驱动设计与用户体验 |
|---|---|---|---|
| AI+开发线 | AI 原生应用开发框架与技术 | AI Agents在研发落地实践 | LLM驱动编程与单测 |
| AI+测试线 | LLM驱动测试分析与设计 | 基于LLM生成测试脚本与数据 | LLM和AI应用的评测 |
| AI+工程线 | AI+DevOps 与工具<br>（LLM 时代的平台工程） | 大模型对齐与安全 | 端侧大模型与云端协同 |
| AI+领域线 | 领域大模型 SFT 与优化 | 知识增强与数据智能 | 大厂专场 |

扫描右侧二维码
查看更多会议详情

**早鸟票**限时抢购中 （截止到9月30日）

¥3680 早鸟票

¥2800 学生票

徐海洋

阿里巴巴通义实验室-高级算法专家

负责通义多模态大模型mPLUG系列工作，包括基础多模态模型mPLUG/mPLUG-2，多模态对话大模型mPLUG-Owl，多模态文档大模型mPLUG-DocOwl，多模态智能体Mobile-Agent等，其中 mPLUG 工作在 VQA 榜单首超人类的成绩。在国际顶级期刊和会议ICML/CVPR/ICCV/ACL/EMNLP/MM/TOIS/IJCAI/AAAI等发表论文40多篇，并担任多个顶级和会议AC/PC/Reviewer。主导参与开源项目mPLUG，X-PLUG，AliceMind，DELTA。

目录
CONTENTS

# PART 01
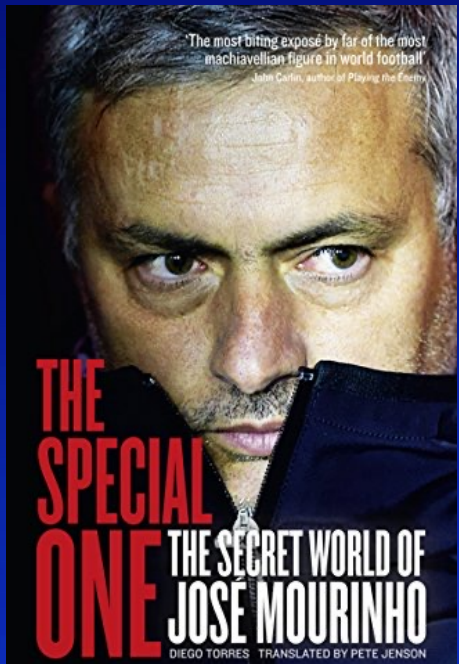# 多模态文档大模型背景

富含文字的图片（简称"文档图片"）
在生活中无处不在!

**Revolutionizing Document AI with Multimodal Document Foundation Models**

## 一个通用的多模态大模型应该具备文档图片理解的能力，例如 GPT4V



**Prompt:**
How many families are earning more than 13000 and owns more than 2 cars?

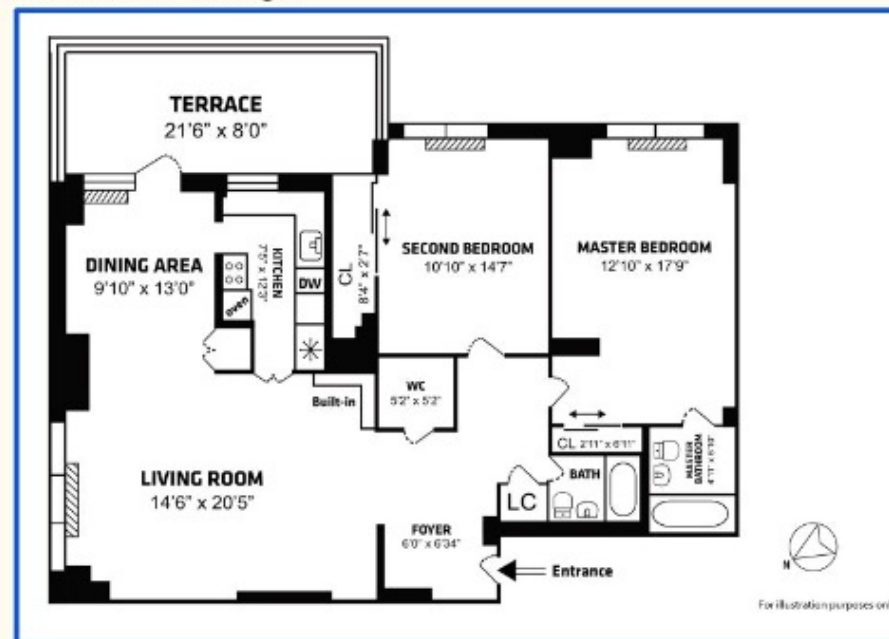| Monthly income (in Rs.) | Vehicles per family | | | |
|---|---|---|---|---|
| | 0 | 1 | 2 | Above 2 |
| Less than 7000 | 10 | 160 | 25 | 0 |
| 7000 – 10000 | 0 | 305 | 27 | 2 |
| 10000 – 13000 | 1 | 535 | 29 | 1 |
| 13000 – 16000 | 2 | 469 | 59 | 25 |
| 16000 or more | 1 | 579 | 82 | 88 |

**GPT-4V:**
From the given table, we can clearly see that there are a total of 113 families that earn more than Rs. 13000 and own more than 2 cars (25 families in the 13000-16000 range and 88 families in the 16000 or more range).



**Prompt:**
Describe the image.

**GPT-4V:**
This image is a floor plan for an apartment or condominium. It shows the layout of the rooms and the dimensions. There is a terrace, dining area, living room, foyer, two bedrooms, a master bathroom, and a guest bathroom. The floor plan is black and white and is a two-dimensional drawing.

## 多模态大模型模型（MLLM）主要关注于一般图片的对话和问答

### mPLUG-Owl



### LLaVA



### InstructBLIP

由于预训练中存在富含文字的图片，他们也初显一定的图片文字理解能力

**MLLM具备实现通用文档图片理解的潜力**

**利用MLLM进行文档图片理解的五大挑战：**

- 文字相关的图片在类型、形状、大小三个方面都十分多样化，已有的MLLM的视觉编码器难以编码这些图片
  2023.7 mPLUG-DocOwl/UReader：无参数的形状适应的切图模块，EMNLP2023

- 文档相关的任务形式很多样，包括信息抽取，问答和自然语言推理等，涉及文字识别和语义理解等不同层次；
  2023.7 mPLUG-DocOwl/UReader：多任务联合学习，EMNLP2023

- 文档图片在结构布局方面十分复杂多样化，其对于文档语义理解十分关键；
  2024.3 mPLUG-DocOwl 1.5：统一的结构学习

- 作为文档图片的一类，Chart问答更考验模型的多步推理和数学计算的能力；
  2024.4 TinyChart-3B：通过Program-of-Thought解决Chart领域多步推理和数学计算问题

- 特殊场景（例如论文理解）下，文档图片的理解需要外部知识的支持；
  2023.12 mPLUG-PaperOwl：基于上下文的论文图表分析，MM2024

# PART 01
# 多模态文档大模型背景

**DocOwl：通过轻量化训练，将已有的MLLM拓展到文档图片领域**

1. 文字相关的图片在类型、形状、大小三个方面都十分多样化，已有的MLLM的视觉编码器难以编码这些图片；

   Ours: 无参数的"形状适应的切图模块"

2. 文档相关的任务形式很多样，包括信息抽取，问答和自然语言推理等，涉及文字识别和语义理解等不同层次；

   Ours: 多任务联合学习+辅助学习任务

Ye J, Hu A, Xu H, et al. **mplug-docowl**: Modularized multimodal large language model for document understanding[J]. arXiv preprint arXiv:2307.02499, 2023.

Ye J, Hu A, Xu H, et al. **UReader**: Universal OCR-free Visually-situated Language Understanding with Multimodal Large Language Model[J]. arXiv preprint arXiv:2310.05126, 2023.

- 模型结构



形状适应的切图模块
1. $Grids = \{1\times1, 1\times2, \dots 1\times20,$

$2\times1, 2\times1, \dots, 2\times10,$

$\dots$

$20\times1\}$
2. 同时考虑形状和大小进行grid选择）

每个子图在原图中第几行第几列

- 模型训练: 训练任务



下游任务

Visual Question Answering (问答)

Information Extraction (信息抽取)

Natural Language Inference (自然语言推理)

Image Captioning (图像描述生成)

辅助任务

Text Reading (文字阅读)

Key Points Generation (要点生成)

Human: What will the presentation aim to encourage a discussion about
AI: It will aim to encourage a discussion about open data ownership rights and the differing social

Human: what is the value mentioned in the CID#:
AI: 736 00 00

Human: juan román riquelme have the same score on july 11 , 2007 november 17 2007, Yes or No?
AI: Yes

Human : Create a compact narrative representing the image presented.
AI: Sign that says Calea Dorobantilor in blue and white

- 模型训练: 训练任务

**辅助任务**

**Text Reading: 增强模型基础的文字识别能力**

- 整个图片的文字序列按照从上到下，从左到右的形式组织，
- 随机选取序列中$\{0, \frac{1}{6}, \frac{2}{6}, ..., \frac{5}{6}\}$的一个位置（概率分布为$\{0.5, 0.1, 0.1, ..., 0.1\}$）作为断点,
- 要求模型根据图片和断点前面的文字，生成断点后面的文字
- 相比所有的样例都读全部文字，这种方法可以让模型更多地关注图片上不同区域的文字

**Key Points Generation：增强模型更高层次的语义理解能力**

- 对于每张图片的question-answer对，使用Vicuna将其转为陈述句的形式，即key point；
- 要求模型根据图片，生成多个key points

- 模型训练: 训练数据

- 实验结果

Table 1: Comparison with ocr-free methods on various types of visually-situated language understanding tasks. 'TSFT' means task-spcific fine-tuning on the downstream dataset. 'underline' means achieving 80% SOTA performance.

| Model | Train Param | TS FT | Doc VQA | Info VQA | Deep Form | KLC | WTQ | TabFact | ChartQA | TextVQA | TextCaps | Visual MRC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Dessurt | 127M | ✓ | 63.2 | - | - | - | - | - | - | - | - | - |
| Donut | 176M | ✓ | 67.5 | 11.6 | 61.6 | 30.0 | 18.8 | 54.6 | 41.8 | 43.5 | 74.4 | 93.91 |
| Pix2Struct$_{base}$ | 282M | ✓ | 72.1 | 38.2 | - | - | - | - | 56.0 | - | 88.0 | - |
| Pix2Struct$_{large}$ | 1.3B | ✓ | **76.6** | 40.0 | - | - | - | - | 58.6 | - | 95.5 | - |
| DocOwl | 86M | × | 65.4 | **42.2** | 49.5 | **32.8** | **29.4** | **67.6** | **59.3** | **57.6** | **118.4** | **221.7** |

结构布局理解对于多模态
文档理解十分关键



(a) What is the assigned response code? W24

(b) which edition has unlimited remote desktop services and virtulization rights? Datacenter

(c) What is the title of the paper in the website? Balance Control in Obese Subjects during Quiet Stance: A state-of-the Art.

(d) What percentage of teenagers from the age group 12-17 didn't use the Internet? 5%

(e) What is the forecast for the increase in customs duty revenue in 2030? 100

(f) What is the percentage of Iraqi dependents citizen? 0.19.

Figure 2: Illustrations of the importance of structure information in Visual Document Understanding on documents (a), tables (b), webpages (c), infographics (d), and charts (e-f).

Hu A, Xu H, et al. mPLUG-DocOwl 1.5: Unified Structure Learning for OCR-free Document Understanding [J]. arXiv preprint arXiv:2403.12895.

DocOwl 1.5: 通过统一的结构学习，显著提升多模态文档理解效果。

**模型结构：视觉语言对齐结构 HReducer**
- 通过卷积更好的保持视觉特征的空间位置关系；
- 降低高分辨率视觉特征的数量；

**模型训练: Unified Structure Learning**
- 结构感知的文字解析
- 多粒度文字识别和定位；

Hu A, Xu H, et al. mPLUG-DocOwl 1.5: Unified Structure Learning for OCR-free Document Understanding [J]. arXiv preprint arXiv:2403.12895.

# mPLUG-DocOwl1.5



Hu A, Xu H, et al. mPLUG-DocOwl 1.5: Unified Structure Learning for OCR-free Document Understanding [J]. arXiv preprint arXiv:2403.12895.

延续DocOwl, 采用Shape-Adaptive Cropping Module处理高分辨率图片

Hu A, Xu H, et al. mPLUG-DocOwl 1.5: Unified Structure Learning for OCR-free Document Understanding [J]. arXiv preprint arXiv:2403.12895.

Twitter. According to the image, there are 560 million and 70 million active users for Twitter and Pinterest. Thus, Twitter has more active users.

Large Language Model  MAM

$T_0$  $T_1$  $T_2$  $T_3$  $T_4$  $T_5$  $T_6$  $X$

<global-img>  <row1-col1>  <row1-col2>  <row1-col3>  <row2-col1>  <row2-col2>  <row2-col3>

$\hat{V}_0$  $\hat{V}_1$  $\hat{V}_2$  $\hat{V}_3$  $\hat{V}_4$  $\hat{V}_5$  $\hat{V}_6$

H-Reducer

$V_0$  $V_1$  $V_2$  $V_3$  $V_4$  $V_5$  $V_6$

Visual Encoder

$I_0$  $I_1$  $I_2$  $I_3$  $I_4$  $I_5$  $I_6$

Shape-Adaptive Cropping Module

Who has more active users, Pinterest or Twitter?

H-Reducer

Convolution (1x4)

$V_i$  $\bar{V}_i$  flatten  FC  $\hat{V}_i$

- 不同于DocOwl, 采用基于卷积的HReducer将视觉特征对齐到语言特征
- 考虑到文字大部分为横向排布，采用1x4的卷积

Hu A, Xu H, et al. mPLUG-DocOwl 1.5: Unified Structure Learning for OCR-free Document Understanding [J]. arXiv preprint arXiv:2403.12895.

- 不同于DocOwl采用学习的位置编码, DocOwl 1.5直接采用特殊的文本token表示当前子图的位置

Hu A, Xu H, et al. mPLUG-DocOwl 1.5: Unified Structure Learning for OCR-free Document Understanding [J]. arXiv preprint arXiv:2403.12895.

• 模型训练



Stage 1: Unified Structure Learning

Stage 2: Multi-task Tuning

Hu A, Xu H, et al. mPLUG-DocOwl 1.5: Unified Structure Learning for OCR-free Document Understanding [J]. arXiv preprint arXiv:2403.12895.

- 模型训练



Figure 4: The illustration of Unified Structure Learning of DocOwl 1.5.

Hu A, Xu H, et al. mPLUG-DocOwl 1.5: Unified Structure Learning for OCR-free Document Understanding [J]. arXiv preprint arXiv:2403.12895.

- 模型训练



| Task | Dataset | Num |
|---|---|---|
| Document Parsing | CCpdf | 938,058 |
| | RVL-CDIP | 159,418 |
| | DUE | 56,809 |
| | VisualMRC | 7,016 |
| Table Parsing | TURL | 200,000 |
| | PubTabNet | 200,000 |
| Chart Parsing | ChartQA | 18,317 |
| | FigureQA | 99,714 |
| | PlotQA | 157,070 |
| | DQA | 200,000 |
| Natural Image Parsing | OCR-CC | 1,000,000 |
| Multi-Grained Text Localization | Word | 114,344 |
| | Phrase | 113,130 |
| | Line | 378,432 |
| | Block | 394,094 |
| Total | | 4,036,402 |

Figure 5: Detailed statistics of DocStruct4M.

Hu A, Xu H, et al. mPLUG-DocOwl 1.5: Unified Structure Learning for OCR-free Document Understanding [J]. arXiv preprint arXiv:2403.12895.

Table 1: The detailed statistics of DocReason25K. The 'Avg Length' refers to the average token length of the answer.

|  | DocVQA | InfoVQA | WTQ | VisualMRC | ChartQA | TextVQA | ALL |
|---|---|---|---|---|---|---|---|
| Image | 1,491 | 1,614 | 850 | 1,927 | 1,252 | 1,612 | 8,746 |
| Sample | 5,119 | 5,421 | 5,994 | 5,263 | 1,827 | 2,253 | 25,877 |
| Avg Length | 79.2 | 95.4 | 77.7 | 103.4 | 106.9 | 88.0 | 89.9 |

已有下游任务只提供简单的词组作为答案，**并不能给出详细的解释**

DocReason25K：基于GPT3.5/GPT4V构建的多模态文档领域关于详细解释的微调数据集

**DocOwl 1.5-Chat**: 联合已有下游任务和DocReason25K进行第二阶段训练。

Hu A, Xu H, et al. mPLUG-DocOwl 1.5: Unified Structure Learning for OCR-free Document Understanding [J]. arXiv preprint arXiv:2403.12895.

Table 2: Different settings of OCR-free Visual Document Understanding models. 'Open' refers to whether all OCR learning data is open-source.

| Model | Init | Resolution | Text | Bbox | OCR Learning Size | OCR Learning Domain | Open |
|---|---|---|---|---|---|---|---|
| Donut [23] | - | 2560x1920 | ✓ | × | 13M | Synthetic, Doc | ✓ |
| Pix2Struct [24] | - | $2^{19}$(shape variable) | ✓ | × | 80M | Web | × |
| QwenVL [3] | - | 448x448 | ✓ | × | 24.8M | Synthetic, Doc, Web | × |
| Monkey [26] | QwenVL [3] | 896x896 | × | × | - | - | - |
| UReader [57] | Owl [58] | 224x224(x20 crops) | ✓ | × | 0.1M | Doc, Table, Chart, Web, Natural | ✓ |
| DocPedia [13] | - | 2560×2560 | ✓ | ✓ | 0.9M | Doc | × |
| CogAgent [46] | CogVLM [51] | 1120×1120 | ✓ | ✓ | 107M | Synthetic, Nature, Doc, Web | × |
| DocOwl 1.5 | Owl2 [59] | 448x448(x9 crops) | ✓ | ✓ | 4M | Doc, Table, Chart, Web, Natural | ✓ |

相比于当前SOTA文档理解大模型：

- DocOwl1.5同时进行了文字、结构以及位置的学习

- DocOwl1.5覆盖了更多的图片类型

- DocOwl1.5只利用开源数据，没有利用私有化数据，易复现

Table 3: Comparison with OCR-free methods on various types of text-rich image understanding tasks. The superscript '*' refers to models separately fine-tuned on each downstream task, rather than generalists. The *underline* means the best performance among models with <10B parameters.

| Model | Size | Doc VQA | Info VQA | Deep Form | KLC | WTQ | Tab Fact | Chart QA | Text VQA | Text Caps | Visual MRC |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Dessurt* | <1B | 63.2 | - | - | - | - | - | - | - | - | - |
| Donut* | <1B | 67.5 | 11.6 | 61.6 | 30.0 | 18.8 | 54.6 | 41.8 | 43.5 | 74.4 | 93.91 |
| Pix2Struct$^*_{base}$ | <1B | 72.1 | 38.2 | - | - | - | - | 56.0 | - | 88.0 | - |
| Pix2Struct$^*_{large}$ | 1.3B | 76.6 | 40.0 | - | - | - | - | 58.6 | - | 95.5 | - |
| DocPeida | 7.0B | 47.1 | 15.2 | - | - | - | - | 46.9 | 60.2 | - | - |
| DocOwl | 7.1B | 62.2 | 38.2 | 42.6 | 30.3 | 26.9 | 60.2 | 57.4 | 52.6 | 111.9 | 188.8 |
| QwenVL | 9.6B | 65.1 | 35.4 | - | - | - | - | 65.7 | 63.8 | - | - |
| UReader | 7.1B | 65.4 | 42.2 | 49.5 | 32.8 | 29.4 | 67.6 | 59.3 | 57.6 | 118.4 | 221.7 |
| Monkey | 9.8B | 66.5 | 36.1 | 40.6 | 32.8 | 25.3 | - | - | 67.6 | 93.2 | - |
| CogAgent | 17.3B | 81.6 | 44.5 | - | - | - | - | 68.4 | **76.1** | - | - |
| DocOwl-1.5 | 8.1B | 81.6 | 50.4 | 68.8 | 37.9 | 39.8 | **80.4** | **70.5** | 68.8 | **132.0** | 239.5 |
| DocOwl-1.5-Chat | 8.1B | **82.2** | **50.7** | **68.8** | **38.7** | **40.6** | 80.2 | 70.2 | 68.6 | 131.6 | **246.4** |

## 样例：文档解析

■ 样例：图表解析

■ 样例：自然图解析



(a a a)

( NIMA RESTAURNT & BAR PARO TOWN Specialist in SHUTNIS DI SM PSTAIRS )

(PEMA DORA GENERAL SHOP)

(a alamy stock photo E7BDR1  www.alamy.com)

Detail the information and textual components depicted in the image.
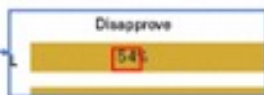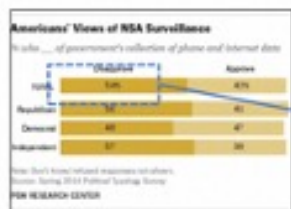
DocOwl 1.5:

the main street in the town <ocr> a a a alamy NIMA RESTAURANT & BAR Special TOURS & DEL PSTAIRS PERIA DORA GENERAL SHOP SONAR ENTERPRISE a a a a a a a alamy photo stock E7BDR1 a www.alamy.com </ocr>

## 样例：多粒度文字识别

■ 样例：多粒度文字定位

■ 样例：详细解释的问答

图表（Chart）理解的主要挑战：

- 问答经常涉及多步数学推理

- 现有最优Chart多模态模型参数太大，推理慢
  （ChartLLama-13B, ChartAssistant-13B）



**College completion by race and ethnicity**

*% ages 25-29 with a bachelor's degree or more, 2014*

Asians 63%
Whites 41
Blacks 22
Hispanics 15

63%
41

Note: Whites, blacks and
include only non-Hispanic
are of any race.
Source: Pew Research Center
tabulations of the March
Population Survey (IPUMS)

PEW RESEARCH CENTER

**Question:** What is the **difference** between
Asians and Whites degree distribution?
**GT Answer:** 22
**ChartLlama:** 21
**TinyChart Direct:** 22

**TinyChart PoT: 22**

```
# Get the values of all
'Values' in the chart, set
to Y
Y=[15, 22, 41, 63]
# Get the difference
between the value of
Asians and Whites, set to
Subtract
Subtract=np.subtract(Y[3],
Y[2])
# Calculate the absolute
value of Subtract, set to
Answer
Answer=np.abs(Subtract)
```

ChartLlama–13B ████ 1.94 it/s
ChartAst–13B ███ 1.47 it/s
TinyChart–3B ██████ 3.14 it/s

(b)

Zhang L, Hu A, Xu H, et al. TinyChart: Efficient Chart Understanding with Visual Token Merging and Program-of-Thoughts Learning[J]. arXiv preprint arXiv:2404.16635, 2024.

Figure 2: Overview of TinyChart.

Zhang L, Hu A, Xu H, et al. TinyChart: Efficient Chart Understanding with Visual Token Merging and Program-of-Thoughts Learning[J]. arXiv preprint arXiv:2404.16635, 2024.

Figure 3: (a) Vision transformer layer with Visual Token Merging. (b) Process of the Visual Token Merging.

Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. 2023. Token Merging: Your ViT But Faster. In The Eleventh International Conference on Learning Representations.

Figure 4: The demonstration of constructing Template-based PoT (upper half) and GPT-based PoT (lower half) in the ChartQA-PoT dataset.

Zhang L, Hu A, Xu H, et al. TinyChart: Efficient Chart Understanding with Visual Token Merging and Program-of-Thoughts Learning[J]. arXiv preprint arXiv:2404.16635, 2024.

**Table 1: Datasets used for training TinyChart. The benchmark datasets consist of basic chart understanding evaluations including QA, summary, and chart-to-table generation. Note that in ablation studies, we only use the benchmark datasets for training due to limited computational resources.**

| Dataset | Benchmark | Samples |
|---|---|---|
| ***Chart question answer*** | | |
| ChartQA [32] | ✓ | 28,299 |
| ChartQA-PoT | ✓ | 140,584 |
| PlotQA [36] | | 157,070 |
| DVQA [19] | | 200,000 |
| OpenCQA [20] | | 5,407 |
| ***Chart-to-text generation*** | | |
| Pew [21] | ✓ | 7,892 |
| Statista [21] | ✓ | 29,589 |
| OpenCQA [20] | | 5,407 |
| Vistext [45] | | 11,171 |
| ChartSumm [42] | | 75,255 |
| Chart2Text-8k [37] | | 7,862 |
| ***Chart-to-table generation*** | | |
| ChartQA [32] | ✓ | 19,373 |
| PlotQA [36] | | 190,720 |
| Chart2Text-8k | | 8,305 |
| DVQA [19] | | 300,000 |
| Statista [21] | | 29,589 |
| ***Chart instruction following*** | | |
| ChartLlama [10] | | 148,398 |
| **Total** | | **1,364,921** |

Zhang L, Hu A, Xu H, et al. TinyChart: Efficient Chart Understanding with Visual Token Merging and Program-of-Thoughts Learning[J]. arXiv preprint arXiv:2404.16635, 2024.

**TinyChart**

Table 2: Main results on chart-related benchmarks. The inference throughput is evaluated on the ChartQA test with a batch size of 1 on V100 32GB.

| Model | #Parameters | Resolution | Inference Throughput | ChartQA | | | Chart-to-Text | Chart-to-Table | OpenCQA |
|---|---|---|---|---|---|---|---|---|---|
| | | | | Aug. | Hum. | Avg. | BLEU4 | RMS$_{F1}$ | BLEU4 |
| *Close source models* | | | | | | | | | |
| GPT-4V [39] | - | - | - | - | - | 78.50 | - | - | - |
| Gemini-Ultra [46] | - | - | - | - | - | 80.80 | - | - | - |
| Qwen-VL-Max [1] | - | - | - | - | - | 79.80 | - | - | - |
| Deplot+Codex [26] | 1.3B+175B | - | - | 91.00 | 67.60 | 79.30 | - | 87.22 | - |
| *Open source models* | | | | | | | | | |
| Llava1.5 [29] | 13B | 336×336 | 1.94 it/s | 72.96 | 37.68 | 55.32 | 7.16 | 48.95 | - |
| Qwen-VL [1] | 9.6B | 448×448 | 1.65 it/s | 78.90 | 44.30 | 61.60 | - | - | - |
| UReader [52] | 7B | 224×224(×20) | 1.67 it/s | 79.42 | 39.12 | 59.30 | - | - | - |
| DocOwl1.5 [15] | 8B | 448×448(×9) | 1.56 it/s | 91.38 | 49.62 | 70.50 | - | - | - |
| ChartInstruct [34] | 7B | - | - | 87.76 | 45.52 | 66.64 | 13.83 | - | 15.59 |
| ChartLlama [10] | 13B | 336×336 | 1.94 it/s | 90.36 | 48.96 | 69.66 | 14.23 | 90.00 | - |
| ChartAst [35] | 13B | 448×448 | 1.47 it/s | **93.90** | 65.90 | 79.90 | 15.50 | 91.60 | 15.50 |
| TinyChart@512 | 3B | 512×512 | **3.65** it/s | 93.60 | 72.16 | 82.88 | **17.93** | 92.93 | 19.62 |
| TinyChart@768 | 3B | 768×768 | 3.14 it/s | 93.86 | **73.34** | **83.60** | 17.18 | **93.78** | **20.39** |

Zhang L, Hu A, Xu H, et al. TinyChart: Efficient Chart Understanding with Visual Token Merging and Program-of-Thoughts Learning[J]. arXiv preprint arXiv:2404.16635, 2024.

Zhang L, Hu A, Xu H, et al. TinyChart: Efficient Chart Understanding with Visual Token Merging and Program-of-Thoughts Learning[J]. arXiv preprint arXiv:2404.16635, 2024.

Zhang L, Hu A, Xu H, et al. TinyChart: Efficient Chart Understanding with Visual Token Merging and Program-of-Thoughts Learning[J]. arXiv preprint arXiv:2404.16635, 2024.

论文图表理解的主要挑战：

1. 需要专业领域知识的辅助

2. 需要具备多图表关联分析能力

A Hu, Y Shi, H Xu,et al. mplug-paperowl: Scientific diagram analysis with the multimodal large language model[J] arXiv preprint arXiv:2311.18248 62023

Figure 2. The pipeline of M-Paper construction and definition of our three tasks.

- 提供**上下文**(context)作为专业知识来源，

* 基于GPT3.5, 构建**"要点"**(outline) 作为图表分析的控制信号

M-Paper: 48k papers, 447k diagrams, 702k training samples

Figure 6: The overall architecture of PaperOwl.

PART 03
mPLUG-DocOwl开源实战

https://github.com/X-PLUG/mPLUG-DocOwl

# ▶ mPLUG-DocOwl Demo

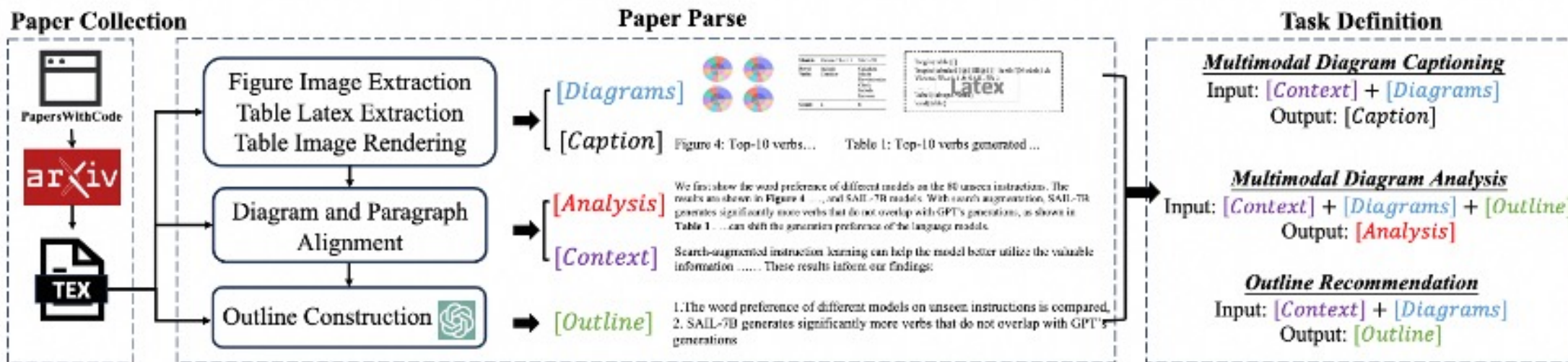https://huggingface.co/spaces/mPLUG/DocOwl

https://modelscope.cn/studios/iic/mPLUG-DocOwl

https://huggingface.co/spaces/mPLUG/TinyChart-3B

# 通义mPLUG多模态体系

AIDD AI+研发数字峰会
AI+ Development Digital summit

ModelScope
Github开源

30+多/单模态任务
SOTA

平台化
业务应用

➤ 图文预训练模型(mPLUG)
提出高效的模态对齐融合新范式 Skip-connected Network 进行跨层/跨语义融合，**CVPR2021 VQA Challenge排名第一**，并以81.26分的成绩首次**超越人类结果**，mPLUG, E2E-VLP相关论文被ACL2021 Oral, EMNLP2022, TIOS接收。

➤ 模块化大一统模型（mPLUG-2）
层次模块化提升单模态和多模态任务水平，**30+多/单模态图文视频任务**，同等数量和模型规模 SOTA效果，在 VideoQA 和 VideoCaption 上 超越 Flamingo、VideoCoca、GITv2 等 超大模型取得 绝对 SOTA，ICML2023接收；

➤ 多模态对话大模型（mPLUG-Owl）
将**模块化轻量化**将视觉表征与大语言模型结合，解锁**大语言模型的多模态能力，多个MLLM榜单第一**，作为多模态LLM通用对比基准；提出mPLUG-Owl2，通过**模态协同同时提升纯文本和多模态性能**，超过 LLaVA1.5，MiniGPT4 等 模型，取得多种模态多种任务 SOTA，CVPR2024接收。

➤ 多模态文档大模型（mPLUG-DocOwl）
模块化轻量化提升各类**文档图表理解**效果，**取得多个Doc任务OCR-Free SOTA**，论文EMNLP2023接收，提出首个论文图文模型mPLUG-PaperOwl，提出文档结构增强的mPLUG-DocOwl1.5，在10个文档理解 benchmark上达到最优效果，5个数据集上提升超过10个点；

➤ 个性化多模态模型（Character-Owl）
多模态个性化理解生成统一大模型，可以进行**轻量化的深度的个性化定制**，具备深度个性化情感多模态能力，CCL Demo演示，mPLUG-Octopus多模态全能助手ACM MM demo，**提名Best Demo**，应用到星尘上的多模态用户发图和角色发图业务，云栖大会发布。

➤ 多模态智能体（mPLUG-Agent）
通过**纯视觉感知**方案操作手机的智能体Mobile-Agent，高效的进行**Self-Planing和Self-Reflection**，不需要依赖XML、系统metadata、系统适配，ICLR 2024 Agent Workshop接收。

mPLUG
多模态深度语言模型体系

端到端图文模型-E2E-VLP
结构化模型StructuralLM
ACL2021 Oral

多模态预训练加速mPLUG-TRIPS
EMNLP2022

多模态大一统模型-mPLUG-2
ICML2023

中文最大视频数据集-Youku-mPLUG
视频对话大模型-mPLUG-video

多模态对话大模型-mPLUG-Owl2
CVPR2024
多模态论文大模型-mPLUG-PaperOwl

多模态文档模型
mPLUG-DocOwl1.5

2021.08    2022.12    2023.07    2024.01

2021.01    2022.05    2023.01    2023.06    2023.11    2024.03

理解生成统一图文模型-mPLUG
CVPR 2021 VQA Challenge
The First Prize
Achieving Human Parity on VQA
TOIS 2022

理解生成统一视频文本模型
mPLUG-HiTeA, TWBERT
ICCV2023

多模态对话大模型-mPLUG-Owl
多模态Denoise-mPLUG-XXM,
ACL2023 Oral

多模态文档大模型-mPLUG-DocOwl
图文预训练加速BUS
ICCV2023

多模态智能体Mobile-Agent

# ▶▶ 通义mPLUG多模态体系

**MobileAgent** `Public`

Mobile-Agent: The Powerful Mobile Device Operation Assistant Family

`android` `agent` `harmony` `ios` `app` `gui` `automation` `mobile` `copilot` `multimodal` `+ 5`

● Python · ⚖ MIT License · ⑂ 179 · ☆ 2.2k · ⊙ 16 · ⇵ 0 · Updated 15 minutes ago

**mPLUG-DocOwl** `Public`

mPLUG-DocOwl: Modularized Multimodal Large Language Model for Document Understanding

`multimodal` `table-understanding` `document-understanding` `mllm` `multimodal-large-language-models` `chart-understanding`

● Python · ⚖ Apache License 2.0 · ⑂ 64 · ☆ 1.1k · ⊙ 27 · ⇵ 1 · Updated 2 weeks ago

**RoleInteract** `Public`

RoleInteract: Evaluating the Social Interaction of Role-Playing Agents

● Python · ⑂ 4 · ☆ 34 · ⊙ 1 · ⇵ 0 · Updated 3 weeks ago

**Multi-LLM-Agent** `Public`

● Python · ⑂ 20 · ☆ 156 · ⊙ 4 · ⇵ 0 · Updated on Apr 23

**mPLUG-Owl** `Public`

mPLUG-Owl & mPLUG-Owl2: Modularized Multimodal Large Language Model

`video` `chatbot` `dialogue` `pytorch` `transformer` `llama` `gpt` `alpaca` `visual-recognition` `multimodal` `+ 10`

● Python · ⚖ MIT License · ⑂ 158 · ☆ 2k · ⊙ 77 · ⇵ 1 · Updated on Apr 6

**SocialBench** `Public`

SocialBench

⚖ Apache License 2.0 · ⑂ 0 · ☆ 1 · ⊙ 0 · ⇵ 0 · Updated on Feb 15

**mPLUG-HalOwl** `Public`

mPLUG-HalOwl: Multimodal Hallucination Evaluation and Mitigating

`benchmark` `contrastive-learning` `hallucinations` `mllm` `multimodal-large-language-models` `multimodal-hallucination`

● Python · ⚖ MIT License · ⑂ 1 · ☆ 59 · ⊙ 3 · ⇵ 0 · Updated on Jan 29

https://github.com/orgs/X-PLUG/repositories

# 通义mPLUG多模态体系



https://modelscope.cn/search?search=mplug

PART 04
总结与展望

- 多模态文档大模型展示了Promising效果，但仍然存在通用性、鲁棒性、效率等问题，仍然还有很多提升的空间；

- 针对高分辨率，多页，探索更高效的视觉编码，图文融合结构；

- 增加更多Domain数据，探索通用的、鲁棒的多模态文档大模型；

- 大小模型协同，能够结合OCR、Chart等专业小模型；



https://github.com/X-PLUG/mPLUG-DocOwl

# 科技生态圈峰会 + 深度研习

## ——1000＋技术团队的共同选择

**KEYLINKing**

### K+峰会

| 🏯 上海站 | 👤 敦煌站 | 🗼 香港站 |
|---|---|---|
| K+峰会 | K+峰会 | K+峰会 |
| **K+全球软件研发行业创新峰会** | **K+思考周®研习社** | **K+思考周®研习社** |
| 时间：2024.06.21-22 | 时间：2024.10.17-19 | 时间：2024.11.10-12 |

K+峰会详情

### AiDD峰会

| 🏯 上海站 | 🏛 北京站 | 🗼 深圳站 |
|---|---|---|
| AiDD峰会 | AiDD峰会 | AiDD峰会 |
| **AI+研发数字峰会** | **AI+研发数字峰会** | **AI+研发数字峰会** |
| 时间：2024.05.17-18 | 时间：2024.08.16-17 | 时间：2024.11.08-09 |

AiDD峰会详情

# AiDD

# 2024 AI+研发数字峰会
## AI+ Development Digital summit

深圳站 11/08-09

## AI 驱动研发变革　促进企业降本增效

### 2024深圳站-议题设置

| AI+产品线 | LLM驱动产品创新 | LLM驱动需求与业务分析 | AI驱动设计与用户体验 |
|---|---|---|---|
| AI+开发线 | AI 原生应用开发框架与技术 | AI Agents在研发落地实践 | LLM驱动编程与单测 |
| AI+测试线 | LLM驱动测试分析与设计 | 基于LLM生成测试脚本与数据 | LLM和AI应用的评测 |
| AI+工程线 | AI+DevOps 与工具（LLM 时代的平台工程） | 大模型对齐与安全 | 端侧大模型与云端协同 |
| AI+领域线 | 领域大模型 SFT 与优化 | 知识增强与数据智能 | 大厂专场 |

扫描右侧二维码
查看更多会议详情

早鸟票限时抢购中（截止到9月30日）

¥3680 早鸟票

¥2800 学生票

# THANKS