# Foundation Models in Medicine Generalist vs Specialist
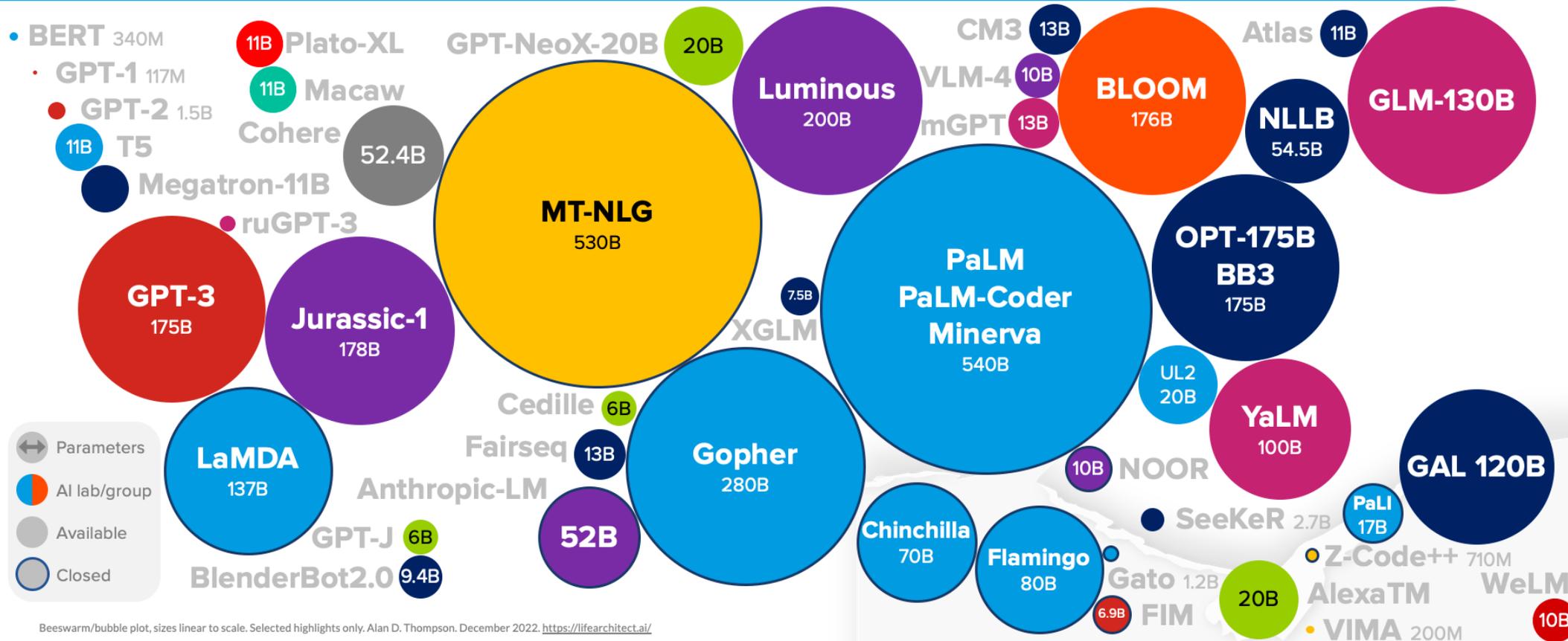
Presented by Shaoting Zhang
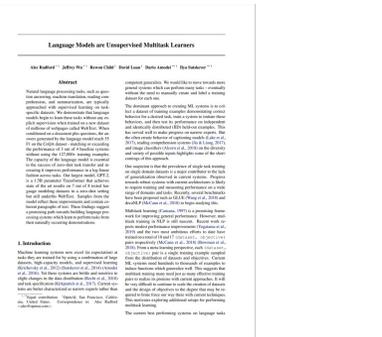
Date：2023.12

# Large-Language Models (LLM)



LANGUAGE MODEL SIZES TO DEC/2022

Beeswarm/bubble plot, sizes linear to scale. Selected highlights only. Alan D. Thompson. December 2022. https://lifearchitect.ai/

LifeArchitect.ai/models

# Large-Language Models (LLM)

**Transformer**
- Proposed the Transformer architecture

**2017/06**

**2018/06**

**GPT**
- Using the decoder of Transformer, a language model is trained on a large amount of unannotated text data, which is then fine-tuned for subtasks

**BERT**
- Using the encoder of Transformer and pre-training on a larger dataset, the performance is much better than GPT

**2018/10**

**2019/02**

**GPT - 2**
- Continuing to use the decoder of Transformer to explore the potential of language models, but the performance is not as impressive as previous models

**GPT - 3**
- Compared to GPT-2, both the data and the model are 100 times larger, with very impressive results

**2020/05**

**2022/11**

**Instruct GPT**
- Through instruction fine-tuning and RLHF (Reinforcement Learning from Human Feedback), the model's generalization ability is enhanced and aligned with human values, making the model controllable and safe
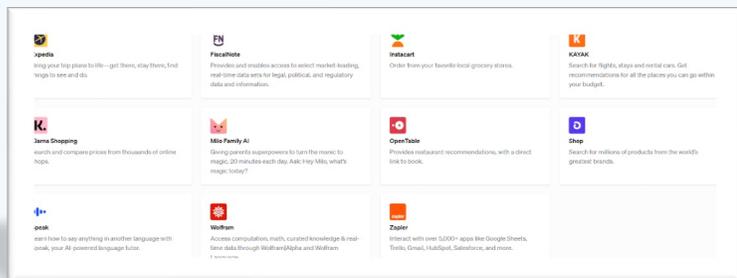
**GPT – 4**
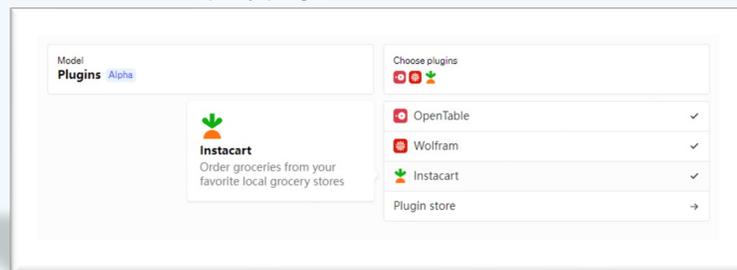- Compared to GPT3.5, it supports image & text input, can analyze images, and provides text output

**2023/03**

# Large-Language Models (LLM)

## Initial support for plugins

➢ The first plugins have been created by Expedia, FiscalNote, Instacart, KAYAK, Klarna, Milo, OpenTable, Shopify, Slack, Speak, Wolfram, and Zapier.



➢ Hosting two plugins, a web browser and code interpreter (Alpha)
➢ Open-sourced the code for a knowledge base retrieval plugin
  • Browsing
  • Code Interpreter
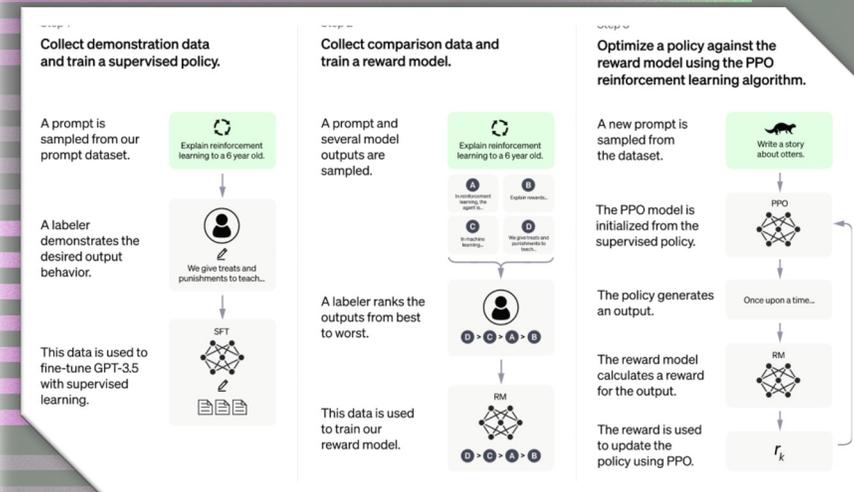  • Retrieval
  • Third-party plugins



**OpenAI**

**ChatGPT**

We've trained a model called ChatGPT which interacts in a conversational way. The dialogue format makes it possible for ChatGPT to answer follow-up questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests.

### The AI model acts as an intelligent API caller

➢ OpenAI plugins connect ChatGPT to third-party applications (Beta)
➢ Plugins enable ChatGPT to interact with APIs defined by developers
➢ Plugins enable ChatGPT to do things like：

  • Retrieve real-time information;
    e.g., sports scores, stock prices, the latest news, etc.
  • Retrieve knowledge-base information;
    e.g., company docs, personal notes, etc.
  • Assist users with actions;
    e.g., booking a flight, ordering food, etc.
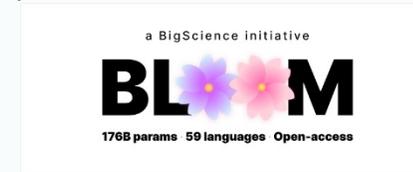
# Open-source Large-Language Models



- **First unified full-stack open-source large-language model system —— INTERN.**
- The InternLM（书生·浦语）200 billion parameter version, InternLM-20B, has been officially open-sourced. The entire toolchain for the development and application of large models has been upgraded, encompassing five major stages: data, pre-training, fine-tuning, deployment, and evaluation.
- InternLM-20B was pre-trained on over 2.3T Tokens containing high-quality English, Chinese, and code data. It exhibits significant improvements in understanding, reasoning, mathematical, and programming abilities
  - ✓ Outstanding overall performance
  - ✓ Strong utility invocation capability
  - ✓ Supports a 16k context length (Through inference extrapolation)
  - ✓ Better value alignment.



| 能力维度 | Llama-13B | Llama2-13B | Baichuan2-13B | InternLM-20B | Llama-33B | Llama-65B | Llama2-70B |
|---|---|---|---|---|---|---|---|
| 语言 | 42.5 | 47 | 47.5 | 55 | 44.6 | 47.1 | 51.6 |
| 知识 | 58.2 | 58.3 | 48.9 | 60.1 | 64 | 66 | 67.7 |
| 理解 | 45.5 | 50.9 | 58.1 | 67.3 | 50.6 | 54.2 | 60.8 |
| 推理 | 42.7 | 43.6 | 44.2 | 54.9 | 46.4 | 49.8 | 55 |
| 学科 | 37.3 | 45.2 | 51.8 | 62.5 | 47.4 | 49.7 | 57.3 |
| 平均成绩 | 43.8 | 47.3 | 49.4 | 59.2 | 48.9 | 51.9 | 57.4 |

本页所列大模型无排名先后
https://github.com/InternLM/InternLM    https://mp.weixin.qq.com/s/
https://mp.weixin.qq.com/s/oTXnvWZJVdoOpFLHngbTYQ



- Parallel training based on RNN.
- Comparable performance with other transformer-based LLMs.





- A series of open-source LLMs free for commercial use（13B & 7B）.
- Based on Transformer, supporting English & Chinses, context window length is 4096.

https://github.com/baichuan-inc/Baichuan-7B    https://github.com/baichuan-inc/Baichuan-13B#readme
https://github.com/BlinkDL/RWKV-LM    https://github.com/THUDM/GLM
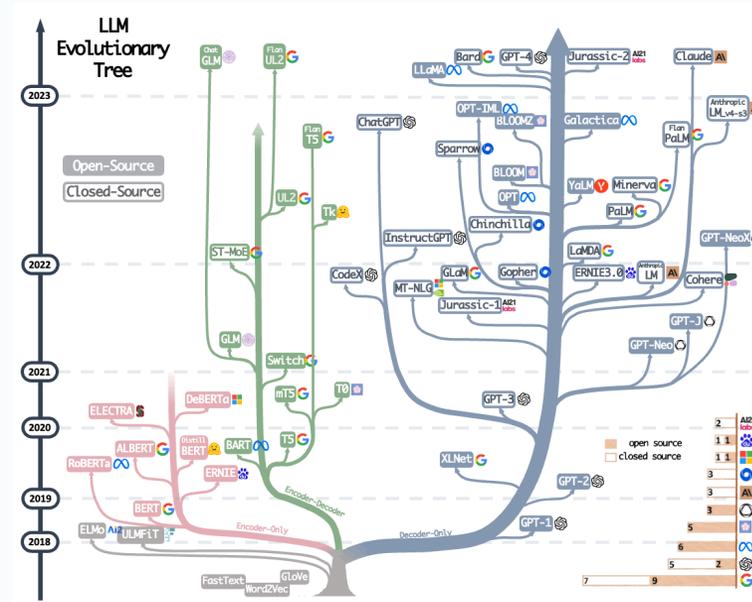https://huggingface.co/RWKV    https://juejin.cn/post/7223305855923044409



- A series of LLMs scaling from 7B to 65B.
- One of the most advanced large-scale foundational language models currently available, and the most famous open-source model.



- 176B LLM which could generate 46 natural languages and 13 programming language texts.
- Using ALiBi position embedding to have more smooth training process and better downstream performance.
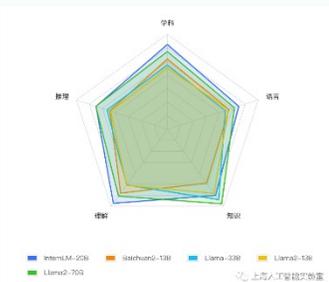- Embedding LayerNorm: Add extra layer norm after embedding layer to improve training stability.



- GLM: A generative LLM, pretraining with auto-regressive blank padding objects and could be fine-tuned in various natural language understanding tasks.
- New ChatGLM2-6B  model: a GLM-based framework, aiming at chinese Q&A and dialogue model, users can deploy in consumer GPUs.

https://github.com/Hannibal046/Awesome-LLM    https://keg.cs.tsinghua.edu.cn/jietang/publications/ChatGLM&Beyond.pdf
https://research.facebook.com/publications/llama-open-and-efficient-foundation-language-models/    https://arxiv.org/pdf/2211.05100.pdf
https://stock.adobe.com/search?k=llama+logo    https://chatglm.cn/blog    https://bigscience.huggingface.co/blog/bloom
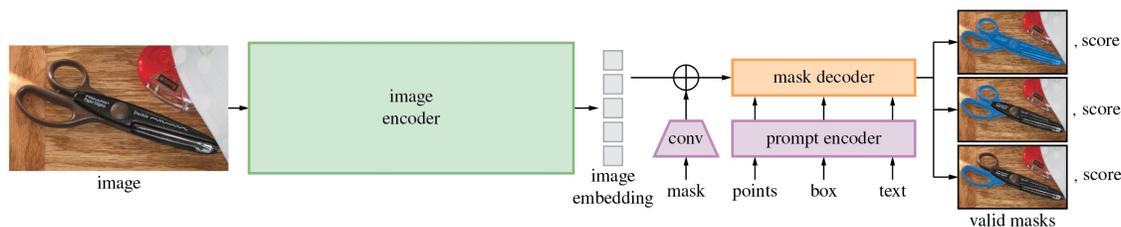
What **task** will enable zero-shot generalization?

What is the corresponding **model** architecture?

What **data** can power this task and model?

**Segment Anything Model (SAM):**

Build a pretrained model with large quantities of parameters in segmentation area, strengthen its generalization by prompt engineering.
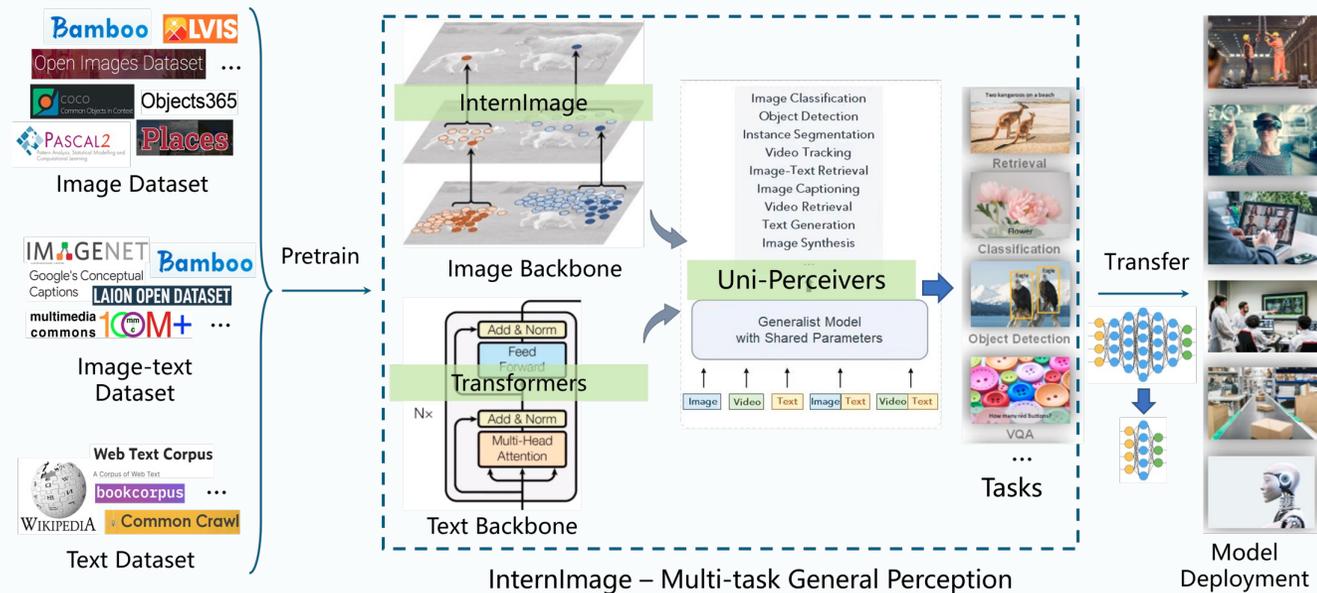
A powerful image encoder computes an image embedding, a prompt encoder embeds prompts, and then the two infor-mation sources are combined in a lightweight mask decoderthat predicts segmentation masks.

Kirillov A, Mintun E, Ravi N, et al. Segment anything[J]. arXiv preprint arXiv:2304.02643, 2023.

https://segment-anything.com/

# Vision Foundation Models

## A new generation of General Vision technology System: INTERN 2.5

"INTERN-2.5" is a powerful multimodal multitask general model jointly released by SenseTime and Shanghai AI Laboratory. It consists of large-scale vision foundation model "InternImage", pre-training method "M3I-Pretraining", generic decoder "Uni-Perceiver" series, and generic encoder for autonomous driving perception "BEVFormer" series.



智慧城市　　自动驾驶　　智慧医疗　　未知领域



Image Dataset

Image-text Dataset

Text Dataset

Pretrain

InternImage

Image Backbone

Transformers

Text Backbone

Nx

Add & Norm
Feed Forward
Add & Norm
Multi-Head Attention

Image Classification
Object Detection
Instance Segmentation
Video Tracking
Image-Text Retrieval
Image Captioning
Video Retrieval
Text Generation
Image Synthesis

Uni-Perceivers

Generalist Model with Shared Parameters

Image Video Text Image Text Video Text

Retrieval
Classification
Object Detection
VQA
... Tasks

Transfer

Model Deployment

InternImage – Multi-task General Perception

Wang W, Dai J, Chen Z, et al. Internimage: Exploring large-scale vision foundation models with deformable convolutions[J]. CVPR 2023.

### Benchmark Table

| 分类 | 图像分类 | | 场景分类 | | | | 长尾分类 | |
|---|---|---|---|---|---|---|---|---|
| | **ImageNet** | | **Places 365** | | **Places 205** | | **iNaturalist 2018** | |
| | BEiT-3 (pub) | ours | ViT-H/14 | ours | MixMIM-L | ours | MAE-H | ours |
| | 89.5 | 90.1 | 60.7 | 61.2 | 69.3 | 71.7 | 86.8 | 92.3 |

| 检测 | 常规物体检测 | | | | 长尾物体检测 | | | | 自动驾驶物体检测 | | 密集物体检测 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **COCO** | **VOC 2007** | **VOC 2012** | **OpenImage** | **LVIS minival** | **LVIS val** | **BDD100K** | | **nuScenes** | | **CrowdHuman** | |
| | ViT-H / ours | Eff-B7 / ours | ATLDET v2 / ours | 冠军性能 / ours | GLIP v2 / ours | Swin-L / ours | PP-YOLOE / ours | | DABNetHD / ours | | Iter-DDETR / ours | |
| | 64.5 / 65.5 | 89.3 / 94.0 | 92.9 / 97.2 | 72.2 / 74.1 | 59.8 / 62.5 | 50.9 / 63.2 | 35.6 / 38.4 | | 62.4 / 64.8 | | 94.1 / 97.2 | |

| 分割 | 语义分割 | | | | | | 街景分割 | | RGBD分割 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | **ADE20K** | | **COCO Stuff-10K** | | **Pascal Context** | | **CityScapes** | | **NYU Depth V2** | |
| | BEiT-3 | ours | ViT-A-L | ours | ViT-A-L | ours | ViT-A-L | ours | CMX-B5 | ours |
| | 62.8 | 62.9 | 54.2 | 59.6 | 68.2 | 70.3 | 85.2 | 86.1 | 56.9 | 69.7 |

| 图文 | 图像描述 | | 微调图文检索 | | | | 零样本图文检索 | |
|---|---|---|---|---|---|---|---|---|
| | **COCO Caption** | | **COCO Caption** | | **Flickr30k** | | **Flickr30k** | |
| | BEiT-3 | ours | BEiT-3 | ours | BEiT-3 | ours | BEiT-3 | ours |
| | 147.6 | 148.2 | 76.0 | 76.4 | 94.2 | 94.8 | 88.2 | 89.1 |

InternImage: Exploring Large-Scale Vision Foundation Models with Deformable Convolutions CVPR2023
Uni-Perceiver v2: A Generalist Model for Large-Scale Vision and Vision-Language Tasks CVPR2023
Towards All-in-one Pre-training via Maximizing Multi-modal Mutual Information CVPR2023
BEVFormer v2: Adapting Modern Image Backbones to Bird's-Eye-View Recognition via Perspective Supervision CVPR 2023
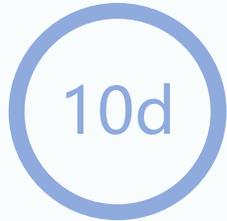
"INTERN-2.5" demonstrated world's best performance on 20 important visual benchmark datasets, covering a wide range of tasks such as classification, detection, and segmentation, making it the top-performing model across multiple domains.

# Foundation Models for Science

## High-precision, Long-term, High-efficiency: AI large model forecasting global weather

Built on multi-modal and multi-task deep learning methods

### Science Large AI Model **FENGWU**

*Global medium-range weather forecast AI Model released by Shanghai AI Laboratory, University of Science and Technology of China, Shanghai Jiao Tong University, Nanjing University of Information Science and Technology, The Institute of Atmospheric Physics, Chinese Academy of Sciences .*

**10d**

Improves the forecast skill and extend the skillful global medium-range weather forecast **out to 10 days** lead **for the first time**.
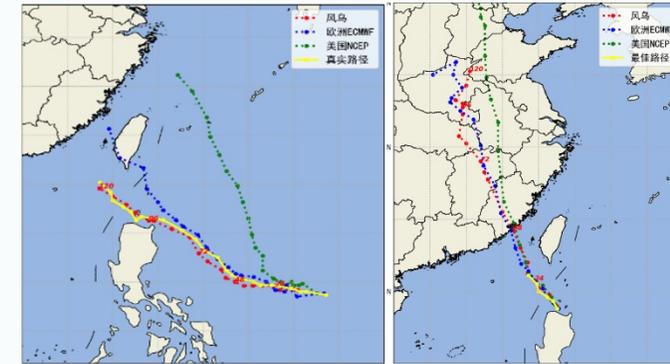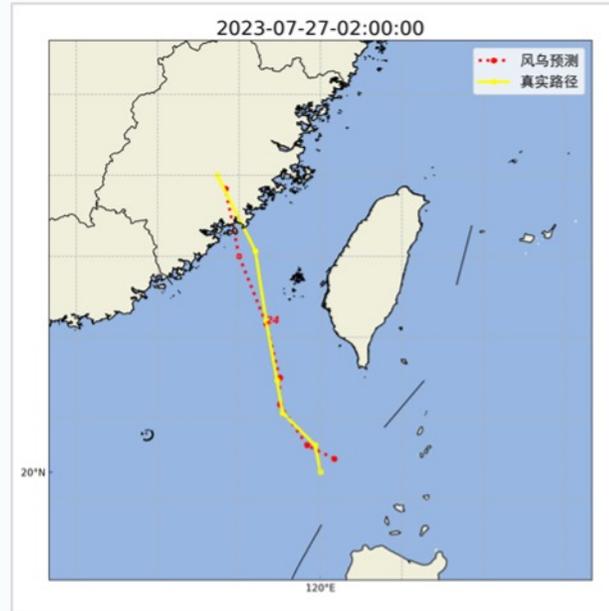
**30s**

Costs **less than 30 seconds** to generate all forecasts in the following 10 days
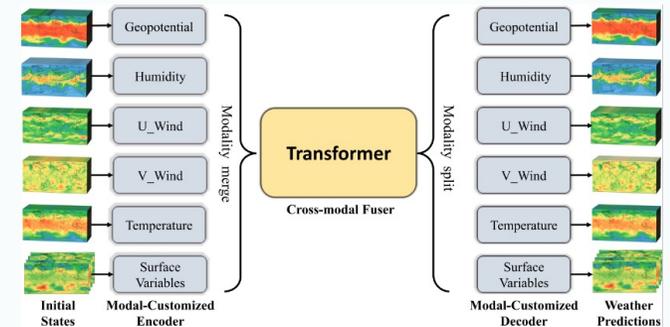
**80%**

Has higher accuracy than **GraphCast** (DeepMind) on **80% reported prediction targets**.

FENGWU: PUSHING THE SKILLFUL GLOBAL MEDIUM-RANGE WEATHER FORECAST BEYOND 10 DAYS LEAD

Chen K, Han T, Gong J, et al. FengWu: Pushing the Skillful Global Medium-range Weather Forecast beyond 10 Days Lead[J]. arXiv preprint arXiv:2304.02948, 2023.

- Recently, with the support of the China Meteorological Administration, **Shanghai AI Laboratory cooperated with the National Meteorological Center and the Shanghai Meteorological Bureau** to complete the real-time trial operation of the Fengwu model deployment. Comparative test evaluation was carried out for the two typhoons "Tali" and "Dusuri" that made landfall this year, and Fengwu showed strong business application potential, providing strong technical support for typhoon forecast and early warning.





Tracking forecasting of Typhoon Doksuri of FengWu, ECMWF and NCEP, figures above are reports from July 22th 8am and July 26th 8am，respectively.



Overview of FengWu's architecture. FengWu first **treats the multiple weather factors as different modalities** and extracts their feature embeddings independently. And then **a transformer-based network is utilized to fuse and passmessages among different modalities.**

K, Singhal, T. Tu, J. Gottweis, et al. Towards Expert-Level Medical Question Answering with Large Language Models. https://arxiv.org/pdf/2305.09617.pdf

https://cloud.google.com/blog/topics/healthcare-life-sciences/sharing-google-med-palm-2-medical-large-language-model

- Med-PaLM is proposed on the basis of Plan-PaLM by prompt engineering method.

- Flan-PaLM achieves SOTA in each MultiMedQA multiple-choice question tasks, its accuracy on MedQA is 67.6%.

- Compared with Plan-PaLM, Med-PaLM has significant improvement in practical problems, its accuracy is comparable with professional doctors, showing its potentials in medical fields.

- Besides Med-PaLM, a new dataset is proposed for medicine LLM evaluation.

Extended Data Fig. 2 | Comparison of SOTA LLMs on MMLU clinical topics. Flan-PaLM achieves state-of-the-art performance on MMLU clinical topics.
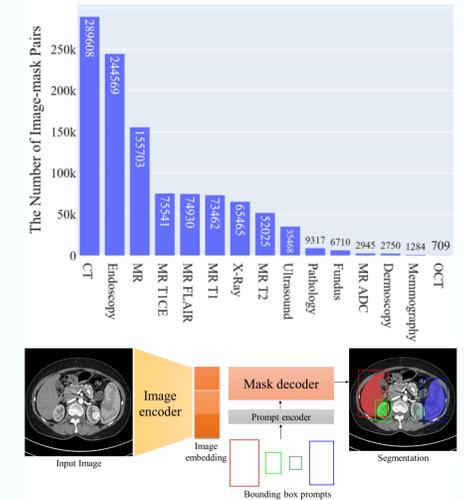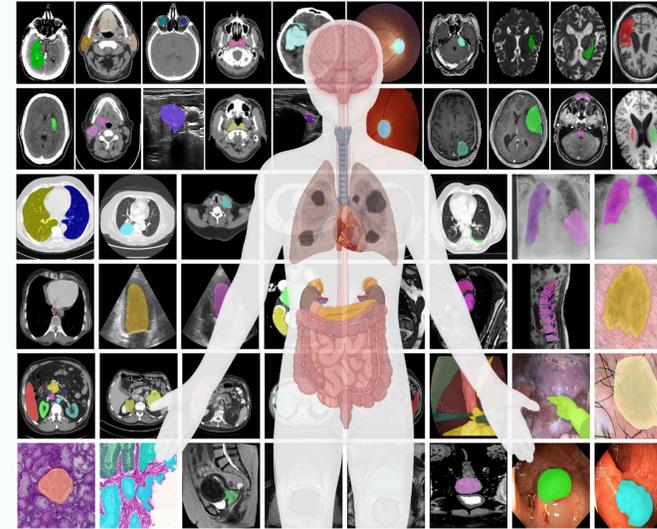
K. Singhal, S. Azizi, A. Karthikesalingam, V. Natarajan, et al. Large language models encode clinical knowledge. Nature. 2023

- MedLAM is the first completely automated medical adaptation of the SAM model, designed to significantly alleviate the annotation workload in the segmentation of medical image.
- Developing MedLAM, a Localize Anything Model for 3D Medical Images
- Coupling this approach with SAM's segmentation capabilities led toan efficient and accurate process for image segmentation.



- Motivated by the remarkable generality of the Segment Anything Model (SAM), MedSAM is the first foundation model for universal medical image segmentation;
- Thoroughly evaluate MedSAM through comprehensive experiments on over 70 internal validation tasks and 40 external validation tasks, spanning a variety of anatomical structures, pathological conditions, and medical imaging modalities;
- MedSAM consistently outperforms the state-of-the-art (SOTA) segmentation foundation model, while achieving performance on par with, or even surpassing specialist models. These results highlight the potential of MedSAM as a powerful tool for medical image segmentation.

W. Lei, X. Wei, X. Zhang, K. Li, S. Zhang. "MedLSAM: Localize and Segment Anything Model for 3D Medical Images" arXiv:2306.14752 2023

Jun Ma, Yuting He, Feifei Li, Lin Han, Chenyu You, Bo Wang. "Segment Anything in Medical Images" arXiv:2304.12306v2 2023

- Protein language model is the transfer application of various language models in the field of biochemistry.

- It can promote the prediction tasks such as protein structure, protein function and generation of new sequences.

- Unsupervised language models can directly utilize massive sequence information in protein databases without relying on manual annotation of sequences.

- Examples: ESM series including ESM-1b, ESM-MSA-1b, ESM-1v, ESM-2 are proposed by Meta.

**ESM-1b**  *family-specific*

**Training:** self-supervised, masked token
**Dataset:** UR50，250 million protein sequence
**Tasks:** Biochemical Properties, Variations, Remote Homology, Secondary Structure and Tertiary Contacts

**ESM-1v**  *zero-shot*

**Training:** self-supervised, masked token
**Dataset:** UR90，protein sequence
**Tasks:** effects of mutations on protein function

**2019**　　**2021**　　**2021**　　**2022**

**ESM-MSA-1b**  *few-shot*

**Training:** self-supervised, masked token
**Dataset:** 4.3 TB，26 million MSAs
**Tasks:** long-range contact、secondary and tertiary structures

**ESM-2、ESMFold**

**Training:** self-supervised, masked token
**Dataset:** UR50，protein sequence
**Tasks:** ESMFold → structure prediction >617 million metagenomic protein sequences

# Foundation Models for Science

- Motivation: the statistical laws of protein sequences imply their structural information → the language model can learn this law through the input sequence.
- Task: filling in missing amino acids in protein sequences → a language model has to understand the underlying structure that creates the patterns in the sequences
- Evidence: As language models of protein sequences are scaled up to 15 billion parameters, an atomic-resolution picture of protein structure emerges in the learned representations.

# Advantage of large models: Emergent abilities

## Emergent abilities

With the advent of large models, the term 'emergence' has become popular. It is typically described as capabilities that do not exist in small-scale models but are present in large-scale models.



Figure 1: **Emergent abilities of large language models**. Language model families display *sharp* and *unpredictable* increases in performance at specific tasks as model scale increases. Emergent abilities [34] have also previously been labeled "abrupt, specific capability scaling" [8], "breakthrough capabilities" [29] and "sharp left turns" [17, 18]. Source: Fig. 2 from [34].

### Are Emergent Abilities of Large Language Models a Mirage?

Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo

Computer Science, Stanford University

## "I am 1.75 meters tall and weigh 50 kilograms. How should I lose weight?"



GPT 3.5

GPT 4

# Advantage of large models: Long-tail problems

Published achievements in digital human medicine, which include 200 anatomical targets of human organs. The data modalities involved include: MRI, Ultrasound, DSA, CT, DR. Over 200+ sections of vascular networks between human organs have been restored, the precision of which surpasses the domain of what human doctors can visually distinguish between arteries and veins. This has reduced the amount of annotated data required for downstream tasks to less than a hundred instances. The model and some finely annotated data are open-source.



**Neurology**

**Cardiovascular**

**Thoracic**

**Skeletal**

**Abdominal**

Precise quantification of brain regions

Ischemic stroke analysis

Intelligent surgical planning for thoracic and abdominal surgery

Lung imaging analysis

Multimodal thoracic and abdominal imaging platform

CT liver analysis

Coronary stenosis analysis

Blood flow reserve evaluation

Fracture localization and morphology analysis

Design of osteotomy guide plate implant

多模态
医学数据

数字人体图谱

打通科室间、
医患间信息壁垒

影像科
智能检出、量化分析

外科
多术智能交互规划

介入科
智能针道规划

Based on a full-chain AI technology stack, we construct a digital human body atlas to achieve efficient and adaptive diagnosis and treatment.

# Advantage of large models: Data security, MaaS

## Model-as-a-Service

| Automated data annotation | Deployment of large model inference | Parallel training of large models | Incremental training of large models | Developer efficiency |
|---|---|---|---|---|
| Intelligent annotation efficiency increased by a hundredfold | Large model inference efficiency improved by 100% | Single cluster with 4000 cards - 500 billion dense parameter model | Incremental fine-tuning cost reduced by 90% | Open model and AI development toolchain |

## Training

- Hardware: 384 pieces of A100, with a single card memory of 80GB, accommodating the batch training of 175 billion parameters

- Cost: $9.6 million (48 servers, $200,000 per server)

- Rental: Single training round ~$3 million (AWS p4d × 2 months)

## Inference

- Hardware: 8 pieces of A100 for single operation, capable of generating about 15-20 words per second (before deployment optimization)

## Estimated Resources

- Training from scratch for a trillion-parameter model: 500-1000 pieces of A100

- Fine-tuning for a trillion-parameter model: 50-100 pieces of A100

### GPT-3.5

| Model | Input | Output |
|---|---|---|
| 4K context | $0.0015 / 1K tokens | $0.002 / 1K tokens |
| 16K context | $0.003 / 1K tokens | $0.004 / 1K tokens |

### GPT-4

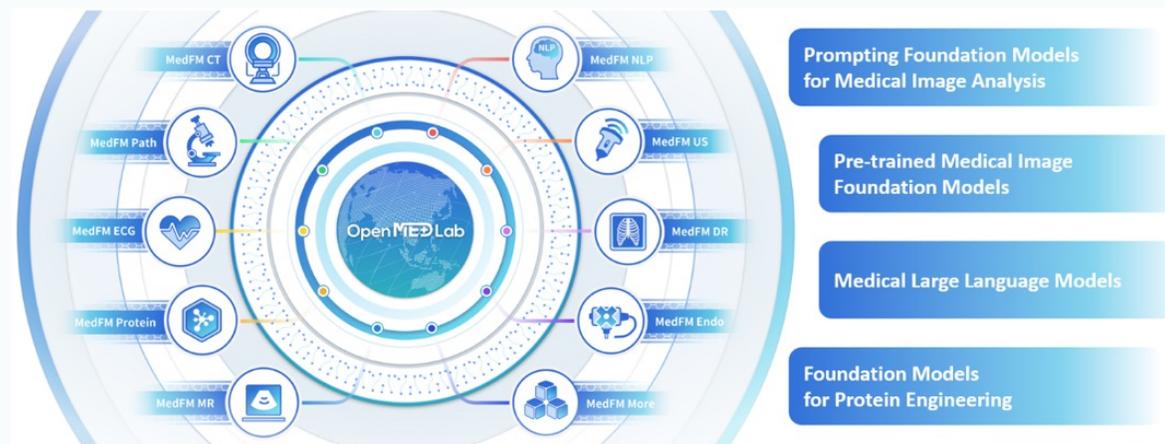| Model | Input | Output |
|---|---|---|
| 8K context | $0.03 / 1K tokens | $0.06 / 1K tokens |
| 32K context | $0.06 / 1K tokens | $0.12 / 1K tokens |

The first multi-level, multi-scenario, high-precision and implementable medical basic model group "OpenMEDLab" has recently been open sourced, integrating the world's top AI research and development capabilities, massive medical data and medical expert knowledge.

The first group of basic models released includes basic models trained on **more than 10 data modalities** such as **medical images**, **medical texts**, **biological information**, and **protein engineering**. Aiming at promoting cross-field, cross-disease, and cross-modal scientific research breakthroughs based on basic medical models, while helping to solve long-tail problems in the medical field and promote the industrial implementation of large medical models.



## World's first open-source platform for medical foundation models "OpenMEDLab"
https://github.com/openmedlab



- The large Chinese medical language model, based on massive medical knowledge and real doctor-patient interaction dialogue data, can provide multi-scenario and multi-round conversation capabilities such as guidance, consultation, health consultation, and assisted decision-making.
- A series of image model groups can achieve cutting-edge research and clinical tasks such as high-precision detection, segmentation, and classification for different modalities such as radiological images, pathological images, endoscopy, and ultrasound images, and efficiently enable solutions to clinical medical imaging-based problems. The model is lightweight and can meet the clinical deployment and convenient application of innovative research.
- In terms of biomedicine, the large-scale protein sequence model is developed. It can help design and synthesize highly stable and bioactive proteins by integrating biomedical expertise, solving the common problems of long time and high cost in the development of new protein drugs.

**The spectrum of foundation models in medical image analysis**



Vision Foundation Model

Medical Foundation Models

Modality-specific Foundation Models

Organ-specific Foundation Models

Task-specific Foundation Models

On the Challenges and Perspectives of Foundation Models for Medical Image Analysis
S. Zhang, D. Metaxas, Medical Image Analysis, 2023, accepted

# Our Related Work: Vision Foundational Model



**Pathology**  **Ultrasound**  **CT**  **MR**  **Electrocardiogram**

**Building global academic influence**

**Medical Image Analysis**

**Special issue on foundational models**
IF 13.8

## Field foundation model 🔥
Applicable to multi-modalities, multi-locations, and multi-diseases

## Small sample weak annotation
Significantly reduce model training costs

### Medical Foundation Models

## High precision ☁
Leading precision in detection, segmentation, and classification

## Efficient transfer 🔗
Rapidly transfer to different tasks

**WAIC 2023 officially released**

**Grand Challenge**

**International challenge**

算法
**Algorithm**

数据
**Data**

算力
**Computational power**

# Our Related Work: Vision Foundational Model

## Grand Challenge
A platform for end-to-end development of machine learning solutions in biomedical imaging.

82,000+ users  •  356 challenges  •  2,750 algorithms

powered by aws

### Three Medical Image Classification Tasks

- Thoracic Disease Screening
- Pathological Tumor Tissue Classification
- Lesion Detection in Colonoscopy Images

https://medfm2023.grand-challenge.org/medfm2023/

| Abnormality | #Sample | Abnormality | #Sample |
|---|---|---|---|
| Cardiomegaly | 1300 | TB | 305 |
| Pleural effusion | 1266 | Pneumothorax | 284 |
| Pneumonia | 898 | Atelectasis | 199 |
| Hilar enlargement | 758 | Emphysema | 159 |
| Nodule | 665 | Calcification | 104 |
| Aortic calcification | 645 | Pulmonary edema | 86 |
| Tortuous aorta | 612 | Increased lung markings | 53 |
| Fibrosis | 532 | Consolidation | 34 |
| Thickened pleura | 497 | Elevated diaphragm | 23 |
| Fracture | 429 | | |

Pleural effusion, Pneumonia, Atelectasis

| Normal Tissue | | Tumor Tissue | |
|---|---|---|---|
| #Sample | 6494 | #Sample | 3515 |
| #Slide | 249 | #Slide | 147 |

| Leison | #Sample | #Study |
|---|---|---|
| Ulcer | 790 | 69 |
| Erosion | 951 | 70 |
| Polyp | 390 | 62 |
| Tumor | 136 | 11 |

Ulcer, Tumor

## MedFMC: Foundation Models for Medical Image Analysis Challenge

This challenge aims to promote technological innovation and explore new learning paradigms in the field of medical image analysis.

Using a small number of data (1-shot, 5-shot, 10-shot) provided by the challenge, transferring the large-scale pre-trained general foundation model to three medical image classification tasks.

(Sponsors)

"A Real-world Dataset and Benchmark For Foundation Model Adaptation in Medical Image Classification"
by X. Wang, D. Wang, L. Wang, M. Li, S. Zhang, et al. Scientific Data. 2023 Accepted

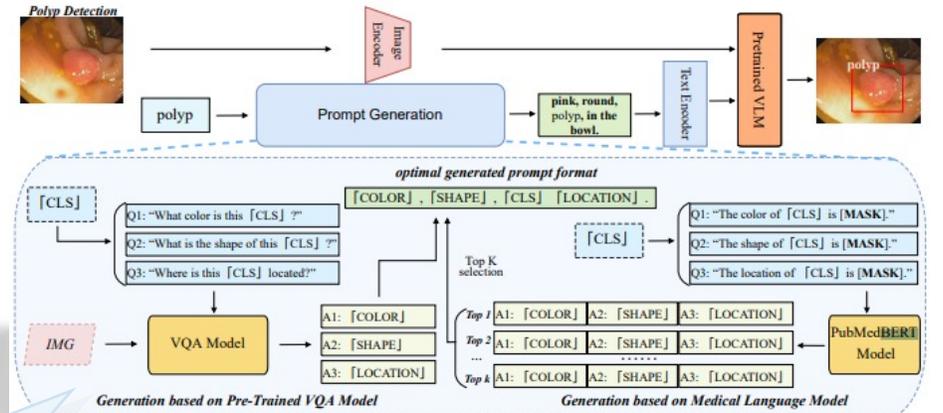## Medical Image Understanding with Pretrained Vision Language Models: A Comprehensive Study

> As one of the pioneering work in applying the visual-language pre-trained model to medical image analysis area, this work first explored how to automatically generate appropriate **language prompts** to guide the pre-trained model identifying the target object in few-shot or even zero-shot situations.

> Using the invariance of descriptive languages such as color and shape in different fields, we can recognize some unseen concepts even in the case of zero-shot tasks. And by inducing a large language model (LLM) to complete adjectives, we successfully automates the generation of these descriptive hints.
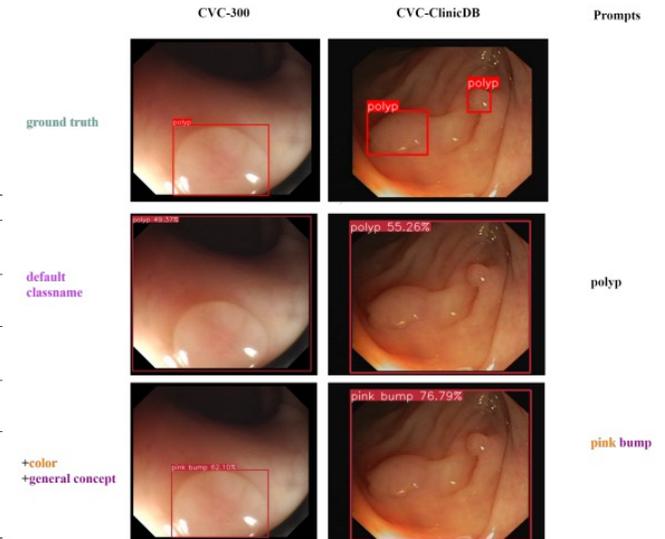
The method in the paper demonstrates the absolute superiority over traditional methods in low-resource scenarios on thirteen medical public datasets spanning multiple modalities.



*Overview of the proposed approach: The optimal medical prompts can be automatically generated with the help of pre-trained VQA model, medical language model, or a hybrid of both.*

Table 2: Our approaches v.s. supervised models on non-radiology datasets (AP%).

| | Method | Backbone | ISIC 2016 | DFUC 2022 | Polyp (×5) | BCCD | CPM-17 | Avg. |
|---|---|---|---|---|---|---|---|---|
| Full Data | Faster RCNN | RN50 | 50.3 | 42.3 | 56.6 | 56.9 | 39.8 | 49.2 |
| | RetinaNet | RN50 | 54.0 | 43.1 | 58.8 | 56.7 | 35.7 | 49.7 |
| | DyHead | Swin-T | 52.9 | 44.2 | 62.9 | 60.1 | 38.8 | 51.8 |
| | GLIP-T(default cls) | Swin-T | 62.4 | 50.3 | 68.1 | 62.5 | 43.9 | 57.4 |
| | Ours (Manual) | Swin-T | 64.1 | 50.3 | 69.4 | 62.2 | 43.4 | 57.9 |
| | Ours (Auto) | Swin-T | 61.6 | 50.1 | 68.8 | 63.1 | 44.2 | 57.6 |
| 100-Shot | Faster RCNN | RN50 | 44.6 | 27.0 | 44.9 | 38.6 | – | 38.8 |
| | RetinaNet | RN50 | 41.7 | 28.4 | 41.7 | 54.3 | – | 41.5 |
| | DyHead | Swin-T | 42.5 | 27.8 | 42.5 | 40.5 | – | 38.3 |
| | GLIP-T(default cls) | Swin-T | 55.9 | 41.4 | 57.6 | 59.8 | – | 53.7 |
| | Ours (Manual) | Swin-T | 58.0 | 43.7 | 60.8 | 60.1 | – | 55.7 |
| | Ours (Auto) | Swin-T | 58.8 | 42.4 | 60.8 | 60.2 | – | 55.6 |
| Zero-Shot | GLIP-T(default cls) | Swin-T | 20.1 | 0.1 | 4.1 | 0.7 | 7.6 | 6.5 |
| | GLIP-L(default cls) | Swin-L | 20.4 | 3.6 | 11.9 | 10.4 | 11.6 | 11.6 |
| | Ours (with MLM) | Swin-T | 25.1 | 24.8 | 38.4 | 24.1 | 20.3 | 26.5 |
| | Ours (with VQA) | Swin-T | 23.5 | 12.9 | 27.1 | 14.3 | 26.2 | 20.8 |
| | Ours (with Hybrid) | Swin-T | 24.5 | 22.5 | 35.1 | 14.3 | 24.8 | 24.2 |
| | Ours (Manual) | Swin-T | 33.3 | 25.9 | 41.3 | 26.9 | 31.4 | 31.8 |



By adding descriptive prompts, the accuracy of model recognition can be greatly improved without additional training

Qin Ziyuan, Huahui Yi, Qicheng Lao and Kang Li. "Medical Image Understanding with Pretrained Vision Language Models: A Comprehensive Study." , ICLR(2023)

## Endo-FM

➢ Constructing a large-scale endoscopic video dataset consisting of over 32K video clips (5M frames), encompassing varying modalities, target organs, and disease types.

➢ Developing a foundation video transformer model capable of capturing both local and global long-range dependencies across spatial and temporal dimensions.

➢ Endo-FM achieves promising performance when applied to various downstream tasks, surpassing state-of-the-art methods by a significant margin.
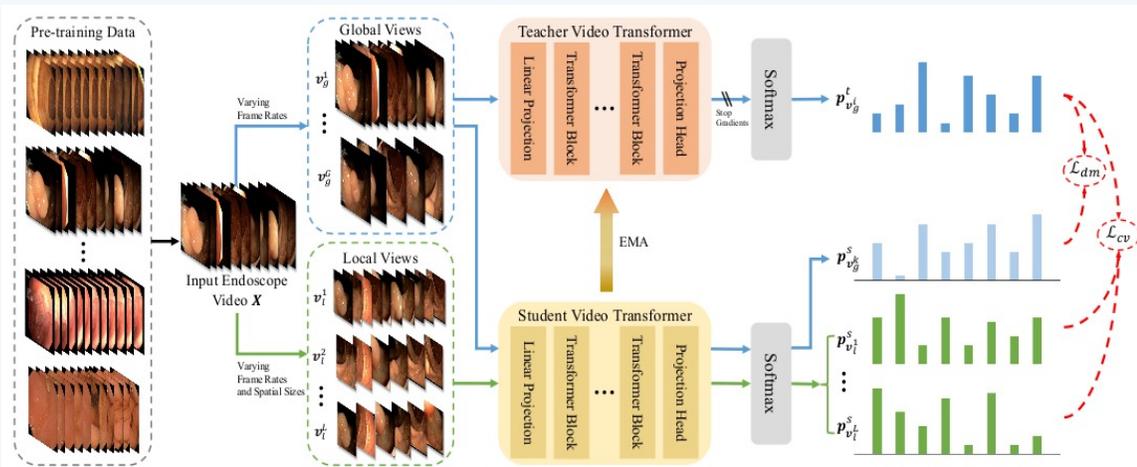
**Table 2.** Comparison with other latest SOTA methods on 3 downstream tasks. We report F1 score for PolypDiag, dice for CVC-12k, and F1 score for KUMC.

| Method | Venue | PolypDiag | CVC-12k | KUMC |
|---|---|---|---|---|
| Scratch (Rand. init.) | | 83.5±1.3 | 53.2±3.2 | 73.5±4.3 |
| TimeSformer [3] | ICML'21 | 84.2±0.8 | 56.3±1.5 | 75.8±2.1 |
| CORP [13] | ICCV'21 | 87.1±0.6 | 68.4±1.1 | 78.2±1.4 |
| FAME [10] | CVPR'22 | 85.4±0.8 | 67.2±1.3 | 76.9±1.2 |
| ProViCo [26] | CVPR'22 | 86.9±0.5 | 69.0±1.5 | 78.6±1.7 |
| Qian et al.[28] | ECCV'22 | 87.6±0.6 | 69.1±1.2 | 78.1±1.9 |
| ST-Adapter [25] | NeurIPS'22 | 84.8±0.7 | 64.3±1.9 | 74.9±2.9 |
| **Endo-FM (Ours)** | | **90.7±0.4** | **73.9±1.2** | **84.1±1.3** |





Example frames of 6 public and ours private endoscope video datasets.

Z. Wang, C. Liu, S. Zhang, Q. Dou. "Foundation Model for Endoscopy Video Analysis via Large-scale Self-supervised Pre-train", MICCAI, 2023.

Article

# A foundation model for generalizable disease detection from retinal images

Yukun Zhou[1,2,3], Mark A. Chia[2,4], Siegfried K. Wagner[2,4], Murat S. Ayhan[1,2,4], Dominic J. Williamson[1,2,4], Robbert R. Struyven[1,2,4], Timing Liu[2], Moucheng Xu[1,3], Mateo G. Lozano[2,5], Peter Woodward-Court[1,2,6], Yuka Kihara[7,8], UK Biobank Eye & Vision Consortium*, Andre Altmann[1,3], Aaron Y. Lee[7,8], Eric J. Topol[9], Alastair K. Denniston[10,11], Daniel C. Alexander[1,12] & Pearse A. Keane[2,4]

- Investigated the development and validation of a foundational model for retinal images, RETFound, with the objective of mitigating the AI model's heavy reliance on massive labeled data and bolstering its generalization capability in disease detection tasks.

- RETFound stands as the inaugural foundational model in the field of ophthalmology, adopting one of the most extensive training datasets and being validating on the most comprehensive set of validation tasks. This open-source model can be tailored to custom downstream tasks, promoting diverse ophthalmological research.

- Trained on 1.6 million unlabeled retinal images by a manner of Self-Supervised Learning and later adapted to disease detection tasks based on explicit labels for the detection of various diseases.

- RETFound outperformed the other pretraining strategies using only 10% of labeled training data and showed consistently high adaptation efficiency, demonstrating the potential of this foundation model in alleviating data shortages and adapting to downstream tasks.

- RETFound provides a generalizable solution to improve model performance and alleviate the annotation workload of experts to enable broad clinical AI applications from retinal imaging.

Y. Zhou, M. Chia, S. Wagner, D. Alexander, P. Keane, et al. "A foundation model for generalizable disease detection from retinal images." Nature, 2023.
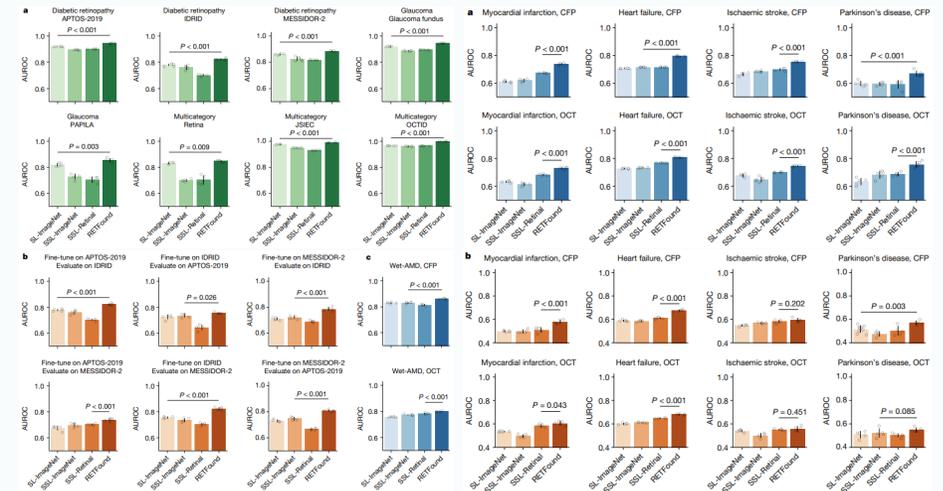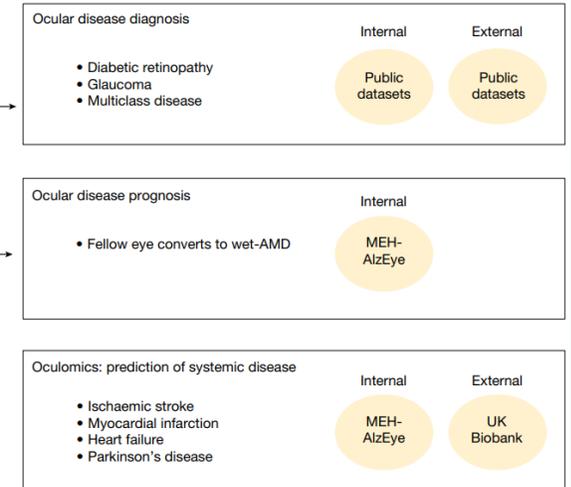https://github.com/rmaphoh/RETFound_MAE   https://new.qq.com/rain/a/20230914A0473700

上海人工智能实验室
Shanghai Artificial Intelligence Laboratory

**Medical Examination Database**

**Medical Quiz 100+**

**Medical Textbook 100+**

**Medical Guidelines 1000+**

**Disease Database with over 4000 Categories (millions of entries)**

**Timely Databases for diseases like COVID-19, Swine Flu, etc**

**Laboratory Test Database**

**Real-world Hospital Consultation Data**

**Real-world Q&A Data (tens of millions of rounds)**

**Real-world Hospital Guide Data**

**Drug Database (millions of entries)**

**Physical Examination Database**

**Adverse Drug Reactions**

**Drug Use Instructions**

## LLM for Medical Foundation Model

- **Large-scale Training:** A model optimized on approximately 4,000,000 instruction datasets from both the Chinese medical domain and general domain.
- **Comprehensive Support for Chinese Medical NLP Tasks:** Including health education, physician examinations, report interpretation, medical record structuring, and simulated diagnostics and treatments.

## MedGPT
### Intelligent consultation display

专业医生的帮助和建议。

您好，感谢您向我咨询。我需要更多的信息才能帮助确定您的症状可能的原因和建议下一步的行动。
首先我想问一下：除了眼睛发黄之外，还有其他不适的症状吗？例如视力模糊、眼痛或头痛等？

最近加班很辛苦，眼睛都黄了。。。

好的，谢谢您提供的额外信息。我的第一个问题是，这情况已经持续了多久了？

眼睛没有不舒服，但是手掌有点黄了

好像有一个多月了

非常抱歉听到您的不适。基于您的描述，这可能是黄疸的症状之一。我建议您尽快预约一次肝功能检查以了解更多情况并排除任何潜在的健康问题。
此外，我也想提醒您注意休息和保持健康的生活方式以帮助缓解症状和促进康复。如果您有任何进一步的问题或需要进一步的医疗咨询，请随时与我联系。

那我应该去哪个科室看病呢

根据您的症状描述，可能需要看内科或者消化科的医生。他们可以对您进行全面的身体检查和相关检测以确定病因并提供适当的治疗方案。
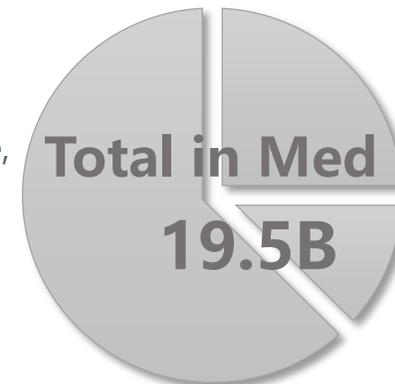同时也可以考虑去眼科做一下视力、验光等方面的相关检查以排除是否存在与视力和眼部疾病相关的因素导致黄疸的出现。

请输入...

## Number of Training Data

Log (x+1)

(bar chart with categories: Self-Supervision, Instruction Tuning, Plugins, RLHF)

Legend: General Domain, Medical Domain

### Self-supervised Pre-training

**Medical Domain Corpus:** Textbook, Guideline, EHR, papers, etc. ~12B tokens
**General Domain Corpus :** Wiki, Webpage, News, etc. ~230B tokens

### Instruction Tuning

**Medical Domain Instruction Tuning Dataset:**
 1 language, 17tasks, 19 types of prompt, 2.7B
**Generic Domain Instruction Tuning Dataset:**
 46 languages, 16 tasks, 16 types of prompt, 5.3M
**Medical Domain QA Dataset:** ~4B tokens
**Generic Domain QA Dataset:** ~1.2B tokens

### Plugins

**General Plugins:** Search Engine, Sci Calculator, ~600M tokens
**Medical Plugins:** Knowledge Retrieval, BMI Calculation, ~109M tokens

### RLHF

**Medical Domain Dialog Dataset:**
 7 scenarios, ~49M tokens
**Generic Domain Dialog Dataset:**
 16 languages, ~72M tokens
**Generic Domain Task Dataset:**
 Chinese & English, 16 tasks, ~35M tokens

**Total in Med 19.5B**

## Medical Domain Supervised Fine-tuning

**Collect demonstration data and train a supervised policy**

- **Instructions:** NER, plugins, etc.

- **Multi-round dialog:** Consultation, triage, etc.

- **QA:** physical examination, drug consultation, etc.

- **Rejection:** refuse to respond non-medical questions

## Reward Model Training

**Collect comparison data and train a reward model**

Design a prompt for self-evaluation

## RL Model Training

**Optimize a policy against the reward model with RL**

Improve the basic PPO framework

Balance the exploration & exploitation



SFT Prompt

RLHF Labeling Tool

RL Models

Standard PPO

Ours

## Customized Knowledge Base

The model can search within a specified knowledge base based on the question posed. It summarizes the content of the related documents retrieved and generates a response to the question.



## Long-range Memory

- **Storing personal information:** No maximum input length restriction, retains long-term health records of individuals.
- **Storing medical information:** Maintains patient information, supports automatic follow-up scenarios.



## Multi-agnt Scheduling

In current complex AI medical diagnostic and treatment processes, there often exist multiple specialized small models, such as data quality control, vessel segmentation, lesion detection, and critical judgment, etc.

Having LLM invoke its own capabilities to complete simple tasks, and coordinating the scheduling of specialized small models for complex tasks, will make complex reasoning across multiple departments and modalities in the diagnostic and treatment phases possible.

**LLM.int8()**

**GPTQ**

**LMDeploy**

**AWQ**

$$s^* = \arg\min_s \mathcal{L}(s), \quad \mathcal{L}(s) = \|Q(W \cdot s)(s^{-1} \cdot X) - WX\|$$

$$s = f(s_X, s_W) = s_X^\alpha \cdot s_W^{-\beta}, \quad \alpha^*, \beta^* = \arg\min_{\alpha,\beta} \mathcal{L}(s_X^\alpha \cdot s_W^{-\beta})$$

| Model Size | Min Hardware Requirements Before Quantization | Min Hardware Requirements After Quantization |
|---|---|---|
| 7B | 1 x 3090 | 1 x 3070 |
| 14B | 1 x A100 | 1 x 3090 |
| 30B | 1 x A100 | 1 x 3090 |
| 100B | 3 x A100 | 1 x A100 |
| 176B | 5 x A100 | 2 x A100 |

# Our Related Work: Large Language Models

**Detailed and accurate "popular science" drug consultation and Q&A capabilities**

**Ability to recommend drugs through light consultation covering common diseases and symptoms**

**Medication consultation LLM**

**Specific training and optimization for authoritative drug data**

**Powerful context-related capabilities and multi-turn dialogue capabilities**

■ **Multi-terminal product application: mobile + PC web**



■ **Medication consultation cases**

Inquiry

Indication

Drug contraindications

Side effect

Drug properties



- The LLM gives medication consultation powerful multi-round dialogue and context-related capabilities.
- Authoritative drug data enables LLM to gain powerful drug understanding and reasoning capabilities
- Simulate conversations between real doctors and pharmacists, eliminating the need to prompt for drug names and related details in each conversation

# Our Related Work: Large Language Models

Based on LLMs and intelligent voice technology, it realizes functions such as batch outbound calls, human-computer interaction, and automatic statistics, helping hospitals complete follow-up of scientific research, specialized disease and caring work, reducing hospital follow-up costs and improving the work efficiency of medical staff.

**Traditional Follow-up pipelines**

Outpatient follow-up, letter-call follow-up, home visit follow-up, entrusted follow-up, telephone and email follow-up

Understand health, morbidity, disease, and death status, and understand satisfaction

**AI for Follow-up**

Follow-up plan

Voice recognition and generation

Follow-up LLM

outbound calls

Follow-up dialogue

| More understandable questions | Analyze non-standard and out-of-range answer content | Understand and connect contextual questions and answers | Thinking and answering out-of-scope questions |

| More humane follow-up interaction | Stronger conversational understanding | More comprehensive information collection |

**Core scientific question: Can a general AI model for protein engineering be constructed to directly design protein products that meet application requirements (end-to-end)?**

1. First, use the sequence and structure information that already exists in nature for pre-training. Our model: ~600 million parameters, billions of protein training data.

2. Then, use the labeled related data set (such as ocean data) for supervised learning.

3. Finally, use a small amount of wet experimental data (dozens) for model tuning.

**Human Corpus Pre-training
(gpt 3.0: 175 billion parameters, hundreds of billions of training data)
Analogy: Chat GPT**

**Standard Dataset Optimization
(supervised learning)**

**Reinforcement Learning for Expert Decision Making**

*Bingxin Zhou, et al. NeurIPS , 2022.* https://arxiv.org/abs/2304.03780

*Zhao，et al., Nature Communications, 13, 7861 (2022)
Li m, et al., Journal of cheminformatics， 15, 12 (2023).*

Professor Hong Liang's Research Group at Shanghai Jiao Tong University

（本页为演示幻灯片，主要为图表内容）

## Crispr Case 12 (gene editing proteins, improve activity)

阳性率55%

IVC

One round of design + experimental iteration, the positive rate of the design is 55%. The editing efficiency of the best mutant is increased by 2 times, achieving the best of its kind.

## T7RNAP (RNA drug polymerase, improves thermostability)

Improve high temperature stability and reduce immunogenic impurity products

Top5 Mutant

It took three and a half months and two rounds of iterations to develop a multi-site mutant with a melting temperature increased by 7 °C. Its stability and activity are comparable to those of the engineering protein developed by the leading company.

**Case 1    Case 2**

**Case 3    Case 4**

## An IVD enzyme (improves thermostability)

单点突变    2点位    3点位    4点位    3点位 4点位 5点位 6点位 8-10点位

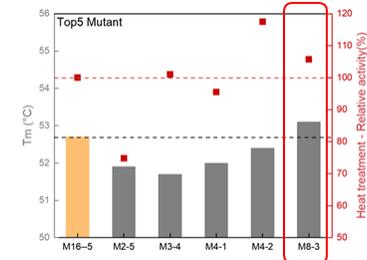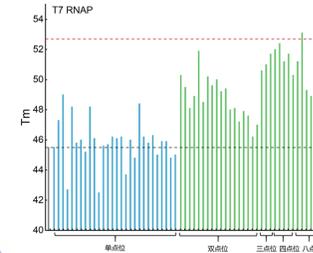Three-round design + wet experiment, the positive rate of the design (30% in the first round), the second round (60%), and the third round (100%). The melting temperature of the best multi-site mutant was 6°C higher than that of the wild type (WT).

## Ago protein (nucleic acid detection/ tumor detection enzyme, improve activity)

1 site   2 sites 3 sites 4 sites 5 sites 6 sites 7 sites 8 sites   9-11 sites

第一轮迭代    第二轮迭代

Two-round design + wet experiment, the positive rate of the design (first round 40%), second round (90%). The optimal mutant has an 8-fold increase in activity compared to the wild type. (WT).

# Our Related Work: Foundation Models for Science

**Institute of Natural Sciences, Shanghai Jiao Tong University**
*JCP, 157 (2022), 144102*
*JCP, 156 (2022), 014114*
*NeurIPS , 2022*

**Innovative Computation and AI Models**

**Shanghai AI Lab**
**Large-scale Natural Language Model + High Computational Power**
*Brie. in Bio., 2022*

**Large-Scale AI Infrastructure and High Computational Power**

General Large-Scale Model for Protein Engine-ering

*J. cheminform.2023.*

Integration of public data, scientific research collaborative data, enterprise cooperative data, totaling: **4.7 billion** unique protein sequence **proprietary database**.
vs the largest public database BFD (**2 billion**)

**Self-built Large Protein Data**

*Nat. Comm., 13, 7861 (2022)*

**Research cooperation**
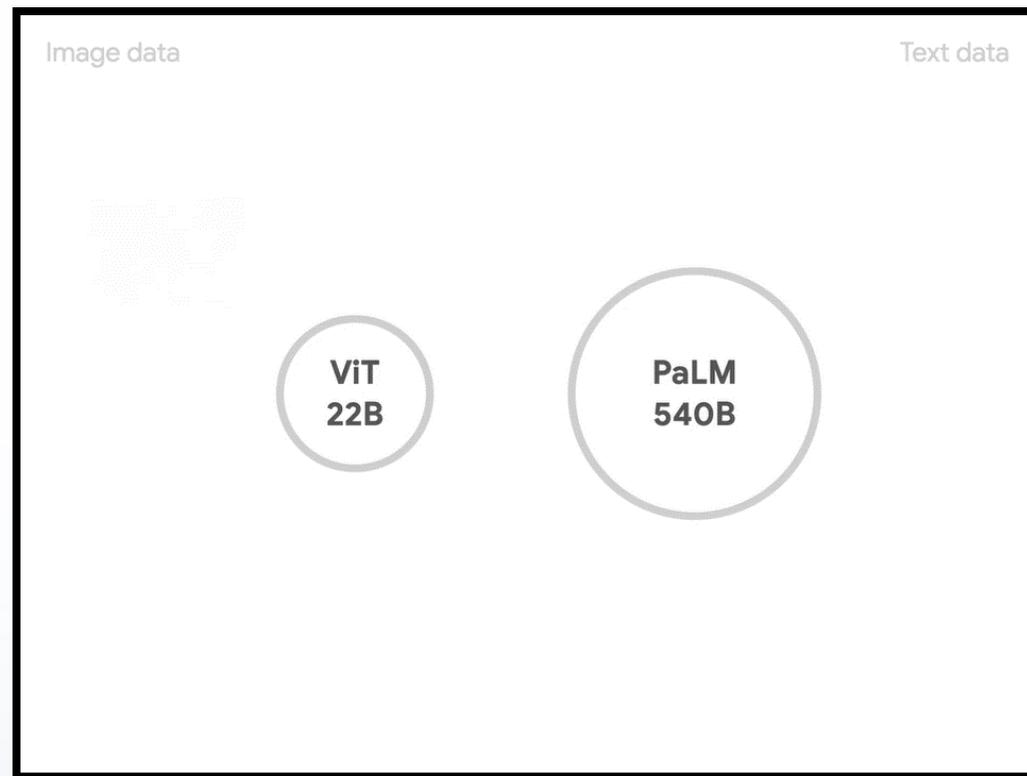1. Stability of a RNA polymerase（**本课题组，阳性)**
2. Working efficiency of a gene editing protein（**上科大免疫所，阳性**)
3. Stability/activity of a nucleic acid detection enzyme（**交大微生物国重室，**
4. Stability of a kidney function test enzyme (**本课题组，阳性**)
5. Activity of non-natural nucleic acid polymerase（**中科院基础医学与肿瘤所，**
6. Activity of artemisinin β-farnesene synthase（**湖北大学酶工程国重室，阴性**

**Enterprise cooperation：**
1. Alkaline resistance of an antibody（**金赛药业**)
2. Substrate selectivity of an engineered enzyme（**瀚海新酶**)
3. Stability of a pharmaceutical enzyme（**康宁杰瑞**)
4. Activity of Nucleic Acid Polymerase（**中元汇吉**)
5. reverse transcriptase （**近岸蛋白**)

| Model type | Large-scale model training in Natural Language Processing | | | |
|---|---|---|---|---|
| **Development Team** | Facebook（USA) | Salesforce Research（USA) | SJTU（Research Group) | |
| **Model Name and Parameter Quantity** | ESM-1V (650 Million Parameters) | ESM2 (150 Million Parameters) | Progen（Xlarge) (6.4 Billion Parameters) | TemPL (650 Million Parameters) |
| **Zero-Shot Model Accuracy** (Correlation coefficient between predictions and laboratory results, for 22 protein TM datasets) | 0.350 | 0.384 | 0.256 | 0.467 |

*Performance Comparison with Internationally Renowned Teams or Companies*

https://arxiv.org/abs/2304.03780

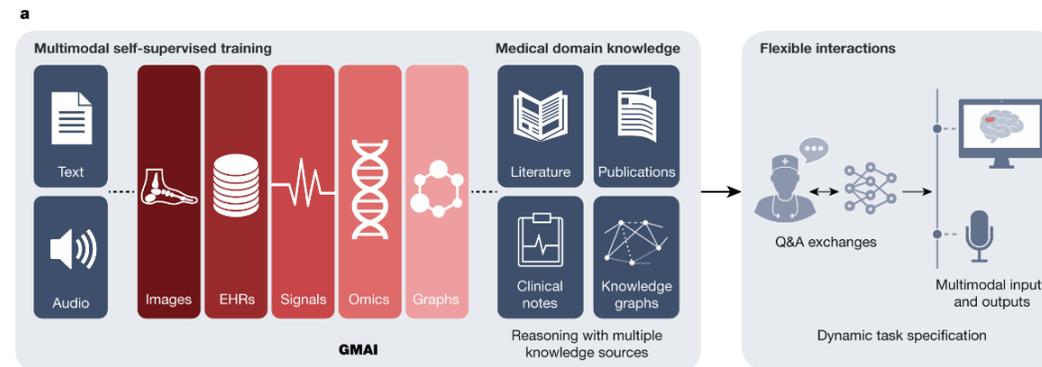Professor Hong Liang's Research Group at Shanghai Jiao Tong University
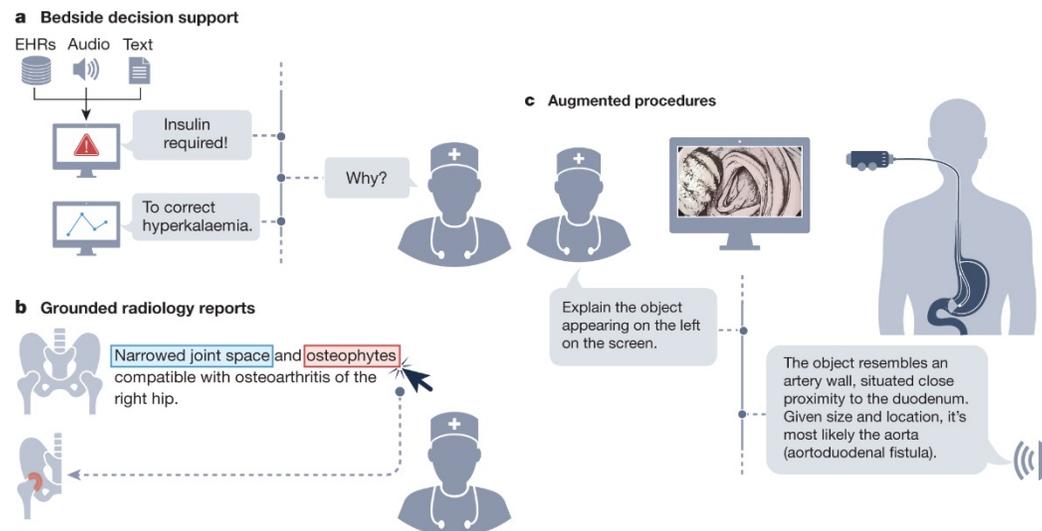
**PaLM-E: A concrete multimodal language model.**

## Foundation models for generalist medical artificial intelligence

- Propose a new paradigm for medical AI, which refer to as generalist medical AI (GMAI)

- Be capable of carrying out a diverse set of tasks using very little or no task-specific labelled data.

- Flexibly interpret different combinations of medical modalities, including data from imaging, electronic health records, laboratory results, genomics, graphs or medical text

- In turn produce expressive outputs such as free-text explanations, spoken recommendations or image annotations that demonstrate advanced medical reasoning abilities

M. Moor, O. Banerjee, Z. Shakeri, et al. Foundation models for generalist medical artificial intelligence. Nature. 2023

**Towards Generalist Biomedical AI**



(1) MultiMedBench, a new multimodal biomedical benchmark, is curated.

(2) MultiMedBench encompasses 14 diverse tasks such as medical question answering, mammography and dermatology image interpretation, radiology report generation and summarization, and genomic variant calling.

(3) Then introducing Med-PaLM M, large multimodal generative model that flexibly encodes and interprets biomedical data including clinical language, imaging, and genomics with thesameset of model weights.

(4) In a side-by-side ranking on 246 retrospective chest X-rays, clinicians express a pairwise preference for Med-PaLM M reports over those produced by radiologists in upto 40.50% of cases, suggesting potential clinical utility.

D. Driess, M. Schaekermann, A. Karthikesalingam, V. Natarajan, et al. Towards Generalist Biomedical AI.
arXiv:2307.14334. 2023

| Task Type | Modality | Dataset | Metric | SOTA | PaLM-E (84B) | Med-PaLM M (Best) |
|---|---|---|---|---|---|---|
| Question Answering | Text | MedQA | Accuracy | 86.50% [61] | 28.83% | 69.68% |
| | | MedMCQA | Accuracy | 72.30% [61] | 33.35% | 62.59% |
| | | PubMedQA | Accuracy | 81.80% [61] | 64.00% | 80.00% |
| Report Summarization | Radiology | MIMIC-III | ROUGE-L | 38.70% [62] | 3.30% | 32.03% |
| | | | BLEU | 16.20% [62] | 0.34% | 15.36% |
| | | | F1-RadGraph | 40.80% [62] | 8.00% | 34.71% |
| Visual Question Answering | Radiology | VQA-RAD | BLEU-1 | 71.03% [63] | 59.19% | 71.27% |
| | | | F1 | N/A | 38.67% | 62.06% |
| | | Slake-VQA | BLEU-1 | 78.60% [64] | 52.65% | 92.7% |
| | | | F1 | 78.10% [64] | 24.53% | 89.28% |
| | Pathology | Path-VQA | BLEU-1 | 70.30% [64] | 54.92% | 72.27% |
| | | | F1 | 58.40% [64] | 29.68% | 62.69% |
| Report Generation | Chest X-ray | MIMIC-CXR | Micro-F1-14 | 44.20% [65] | 15.40% | 53.56% |
| | | | Macro-F1-14 | 30.70% [65] | 10.11% | 39.83% |
| | | | Micro-F1-5 | 56.70% [66] | 5.51% | 57.88% |
| | | | Macro-F1-5 | N/A | 4.85% | 51.60% |
| | | | F1-RadGraph | 24.40% [14] | 11.66% | 26.71% |
| | | | BLEU-1 | 39.48% [65] | 19.86% | 32.31% |
| | | | BLEU-4 | 13.30% [66] | 4.60% | 11.50% |
| | | | ROUGE-L | 29.60% [67] | 16.53% | 27.49% |
| | | | CIDEr-D | 49.50% [68] | 3.50% | 26.17% |
| Image Classification | Chest X-ray | MIMIC-CXR (5 conditions) | Macro-AUC | 81.27% [69] | 51.48% | 79.09% |
| | | | Macro-F1 | N/A | 7.83% | 41.57% |
| | Dermatology | PAD-UFES-20 | Macro-AUC | N/A | 63.37% | 97.27% |
| | | | Macro-F1 | N/A | 1.38% | 84.32% |
| | Mammography | VinDr-Mammo | Macro-AUC | 64.50% [49] | 51.49% | 71.76% |
| | | | Macro-F1 | N/A | 16.06% | 35.70% |
| | | CBIS-DDSM (mass) | Macro-AUC | N/A | 47.75% | 73.31% |
| | | | Macro-F1 | N/A | 7.77% | 51.12% |
| | | CBIS-DDSM (calcification) | Macro-AUC | N/A | 40.67% | 82.22% |
| | | | Macro-F1 | 70.71% [70] | 11.37% | 67.86% |
| Genomics (Variant Calling) | | PrecisionFDA (Truth Challenge V2) | Indel-F1 | 99.40% [71] | 53.01% | 97.04% |
| | | | SNP-F1 | 99.70% [71] | 52.84% | 99.35% |

# About us

## THE LANCET Digital Health

## nature machine intelligence

## nature communications

## scientific data

| 论文标题 | 期刊名称 |
|---|---|
| Spatially aware graph neural networks and cross-level molecular profile prediction in colon cancer histopathology: a retrospective multi-cohort study | The Lancet Digital Health |
| Predicting unseen antibodies' neutralizability via adaptive graph neural networks | Nature Machine Intelligence |
| Mining Multi-Center Heterogeneous Medical Data with Distributed Synthetic Learning | Nature Communications |
| Genetic mutation and biological pathway prediction based on whole slide images in breast carcinoma using deep learning | npj Precision Oncology |
| MedFMC: A Real-world Dataset and Benchmark For Foundation Model Adaptation in Medical Image Classification | Scientific Data |
| A Large-scale Synthetic Pathological Dataset for Deep Learning-enabled Segmentation of Breast Cancer | Scientific Data |

### OpenMED Lab

https://github.com/openmedlab

**WCH, BUPT and UESTC**

Kang Li    Guotai Wang    Qicheng Lao

**Guangzhou & Zhejiang Lab**

Yixue Li    Feng Gao    Wentao Zhu

**SJTU**

Dequan Wang    Xiaofan Zhang    Liang Hong

**CHUK, FDU, and CPII**

Qi Dou    Lijian Xu    Yi Guo

**Shanghai AI Lab**

Junjun He    Xiaosong Wang

**Teams of PIs**

**Overseas**

Mu Zhou    Dimitris Metaxas    Yukun Zhou    Xiaoxiao Li

# Q & A

https://www.shlab.org.cn

Shanghai Artificial Intelligence Laboratory