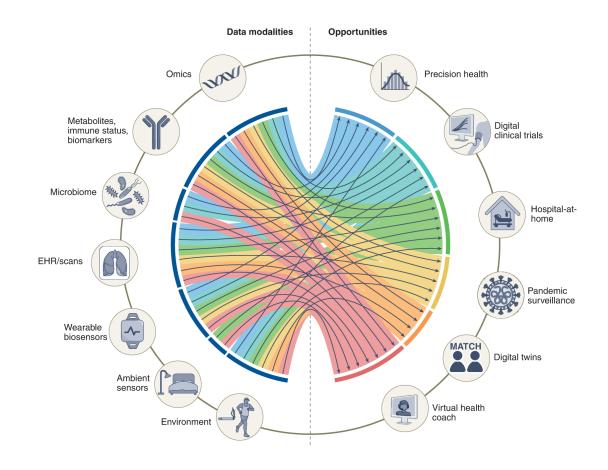




NLP is key to unlocking the future of medicine from clinical text

# Multi-modal is key to unlocking new opportunity

- Health is not limited to a single data modality
- Health is not limited to a single setting
- Holistic models build understanding of health



## Does it work?

# Is it safe?





Real-world data can provide crucial evidence

Human-in-the-loop helps mitigate risk & continuously improve

# Rethinking AI for Health

#### Sparks of Artificial General Intelligence: Early experiments with GPT-4

Sébastien Bubeck Varun Chandrasekaran Ronen Eldan Johannes Gehrke Erie Horvitz Ece Kamar Peter Lee Yin Tat Lee Yuanzhi Li Scott Luundberg Harsha Nori Hamid Palangi Marco Tulio Ribeiro Yi Zhang

Microsoft Research

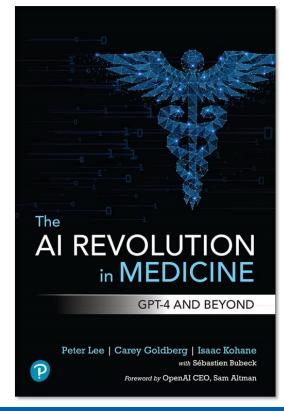
#### Abstract

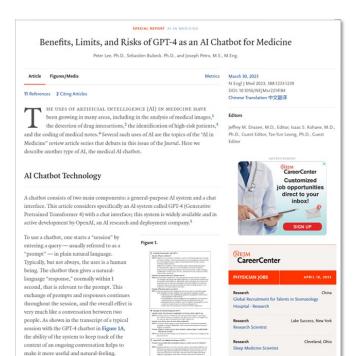
Artificial intelligence (All researchers have been developing and refining large longuage models (LLMs) that exhibit remarkable capabilities acros a variety of domains and tasks, challenging our understanding of learning and cognition. The latest model developed by OpenAI, GPT-4 (Ope-23), was trained using an unprecedented scole of compute and data. In this paper, we report on our investigation of an early version of GPT-4, when it was still in active development by OpenAI. We contend that (this early version of) GPT-4 is part of a new obsort of LLMs (along with ChatGPT and Google's PalM for example) that exhibit more general intelligence than previous AI models. We discuss the rising capabilities and implications of these models. We demonstrate that, beyond its mastery of language, GPT-4 can solve novel and difficult tasks that span mathematics, coding, vision, medicine, law, psychology and more, without needing any special prompting. Moreover, in all of these tasks, GPT-4 per performance is strikingly close to luman-level performance, and often veatly surpasses prior models such as ChatGPT. Given the breatht and depth of GPT-4 is capabilities, we believe the GPT-4 is capabilities, we believe of GPT-4 is predicted and the complete version of AGI, including the possible need for advancing towards deeper and more comprehensive versions of AGI, including the possible need for praviaging new paradigm that moves beyond next-word prediction. We conclude with reflections on societal influences of the recent technological leap and future research directions.

#### Contents

arXiv:2303.12712v5

1	Inti	roduction																								
	1.1	Our approac	h to stud	lying G	PT	-4%	s in	tel	lig	ene	œ															
	1.2	Organization	of our d	lemonst	rat	ion																				
2	Mu	ltimodal and	interd	isciplin	ar	y c	on	аре	osi	tic	n															13
	2.1	Integrative a	bility																							1
	2.2	Vision																								1
		2.2.1 Imag	e generat	ion bey	one	i n	en	ori	iza	tio	n															1
		2.2.2 Imag	e generat	ion foll	owi	ng	det	tail	led	in	sti	uc	tic	ns	(1	la	E	al	l-E	(3						1
		2.2.3 Possi	ble applie	cation i	n sl	ket	ch	ger	ıer	ati	on															1
	2.3																									
	Coc	ling																								2
	3.1	From instru	ctions to	code .																						2
		3.1.1 Codi	ng challer	nges .																						2
			world see																							
	3.2	Understandi	ng existir	ng code																						2





The chatbots in use today are sensitive to the

This aspect of chatbots has given rise to a concept of "prompt engineering," which is

both an art and a science. Although future AI systems are likely to be far less sensitive to the precise language used in a prompt, at present, prompts need to be developed.

form and choice of wording of the prompt

Las Vegas, Nevada

### Use large language models to promote equity

Emma Pierson<sup>1, 2, 14</sup>, Divya Shanmugam<sup>3, 14</sup>, Rajiv Movva<sup>1, 14</sup>, Jon Kleinberg<sup>4, 14</sup>, Monica Agrawal<sup>5</sup>, Mark Dredze<sup>6</sup>, Kadija Ferryman<sup>6</sup>, Judy Wawira Gichoya<sup>7</sup>, Dan Jurafsky<sup>8</sup>, Pang Wei Koh<sup>9</sup>, Karen Levy<sup>4</sup>, Sendhil Mullainathan<sup>10</sup>, Ziad Obermeyer<sup>11</sup>, Harini Suresh<sup>12</sup>, Keyon Vafa<sup>13</sup>

85% of LLM papers studying equity impacts focus on equity harms

- LLMs can improve detection of bias
- LLMs can create structured datasets of equityrelevant information
- LLMs can improve equity of access
- LLMs can improve equity in matching systems

## Insight Consumer Pharma, Payor, Regulator



US: Less than 3% cancer patients enroll in trials 40% cancer trial failures due to insufficient patients New drug costs \$2-10 billion and takes 10+ years



### Insight Consumer

Pharma, Payor, Regulator









Provider, EHR Vendor

⚠ We're building a better ClinicalTrials.gov. Check it	out and tell us what you	u think!				
VIII) U.S. National Library of Medicine  Clinical Trials.gov	Find Studies ▼	About Studies ▼	Submit Studies ▼	Resources ▼	About Site ▼	PRS Login
Home > Search Results > Study Record Detail						Save this study
Cabozantinib in High Grade Neuroendocrine N	leoplasms					
The safety and scientific validity of this study is the responsibility of the study spinvestigators. Listing a study does not mean it has been evaluated by the U.S. For Government. Know the risks and potential benefits of clinical studies and talk to provider before participating. Read our disclaimer for details.		Recruitment St.	osted 1 : December	ı		
			View this study	y on Beta.Clinical	Trials.gov	
Sponsor: Washington University School of Medicine  Collaborator: Exelixis Information provided by (Responsible Party):						
Washington University School of Medicine  Study Details Tahular View No Results Posted Disclaimer 2 How to F	Read a Study Record					
Study Details Tabular View No Results Posted Disclaimer 2 How to F	neau a Study necofd					Go to ▼
Brief Summary:  High grade neuroendocrine neoplasm patients are treated with platinum doublets such as carbordisease succumb to it within a year. There is no extensive literature or consensus on second-or			_			

# LLM: Universal Structuring

### Instruction

Your task is to convert the given clinical trial criteria input into a valid logic formula. Your response should be written in the language of propositional logic and should accurately capture the meaning of the input. Please note that your response should also include any necessary logical connectives, such as "and", "or", or "not". Please keep in mind that your response should be flexible enough to allow for various relevant and creative solutions. You should also focus on providing an accurate and well-structured solution that can be easily understood by others.

#### Input:

"

Histologically confirmed metastatic colorectal adenocarcinoma with mutant APC, TP53 and KRAS genes as determined by the local CLIA-certified laboratory are eligible. All RAS mutations are allowed (KRAS, NRAS, HRAS). Patients with wild type KRAS, APC or TP53 are ineligible.

### Example

### Output:

metastatic AND colorectal adenocarcinoma AND (APC mutation AND TP53 mutation AND KRAS mutation) AND NOT (NOT KRAS mutation OR NOT APC mutation OR NOT TP53 mutation)

Input:

-Histologically or cytologically confirmed high-grade neuroendocrine tumor that has progressed on first line therapy, excluding small cell lung cancer (SCLC). High grade includes any neuroendocrine neoplasm with a Ki-67 of >=20% or with mitotic count of more than 20 mitoses per high power field or any poorly differentiated neoplasm or any neoplasm lacking these that is deemed high grade by pathology consensus, based on other markers (necrosis or IHC demonstrating p53 or RB mutation).

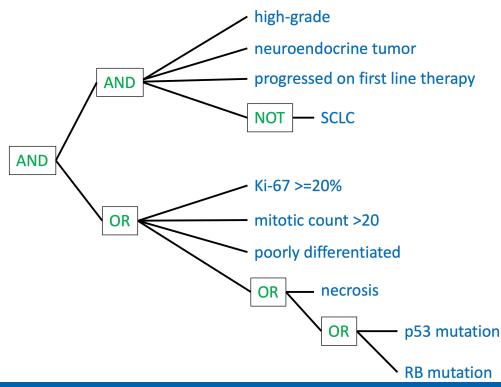
### Input

Output:

# LLM: Universal Structuring

Output

(high-grade AND neuroendocrine tumor AND progressed on first line therapy AND NOT SCLC) AND (Ki-67 >= 20% OR mitotic count >20 OR poorly differentiated OR (necrosis OR (p53 mutation OR RB mutation)))

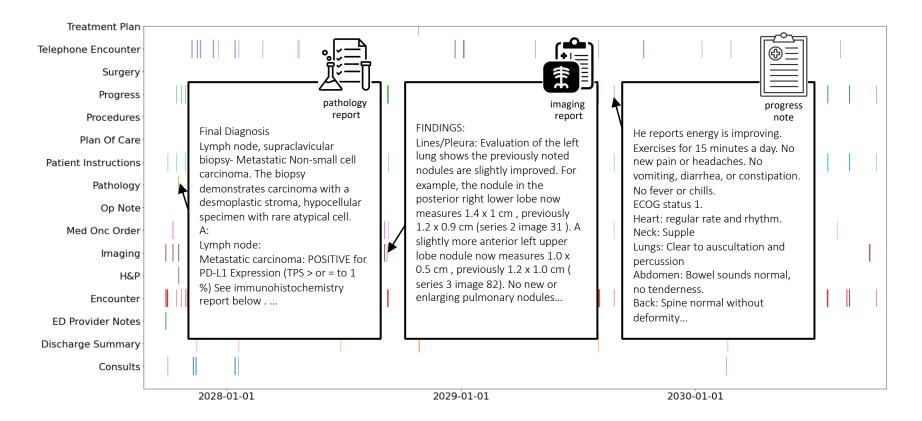


# LLM: Universal Structuring

	His	stology		Bio	marker	
	Precision	Recall	F1	Precision	Recall	F1
GNormPlus	_	_	-	6.8	19.6	10.2
SciSpaCy	34.2	70.2	46.0	58.3	6.9	12.3
Criteria2Query	29.6	40.2	32.8	68.3	27.5	39.2
GPT-3.5 (zero-shot)	35.1	31.6	34.2	61.2	29.4	39.7
GPT-4 (zero-shot)	62.1	69.0	65.4	75.3	59.8	66.7
GPT-4 (3-shot)	57.8	73.7	64.8	72.5	72.5	72.5

Wong et al. "Scaling Clinical Trial Matching Using Large Language Model: A Case Study in Oncology", MLHC 2023.

# EMR: Cancer Patient Journey



# OncoBERT: Oncology RWE



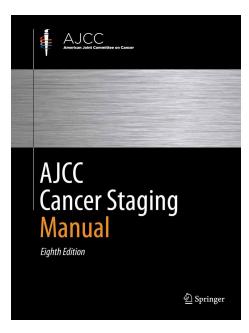
	Tumor Site	Histology	Clinical T	N	M	Pathological T	N	M
Ontology	19.4	19.2	-	-	-	-	-	-
BOW	62.8	76.6	70.4	96.6	98.4	72.1	90.7	98.9
OncoGloVe + CNN	72.0	84.4	74.2	96.5	98.6	83.9	93.1	98.5
OncoGloVe + HAN/GRU	74.0	85.9	76.2	97.1	98.7	86.4	94.2	98.5
BERT + HAN/GRU	75.1	86.2	77.0	96.6	98.4	86.4	94.4	98.2
PubMedBERT + HAN/GRU (ours)	76.7	87.2	79.3	97.2	98.7	87.2	95.2	98.6
OncoBERT + HAN/GRU (ours)	77.1	87.6	81.4	97.5	99.0	87.6	95.5	98.9

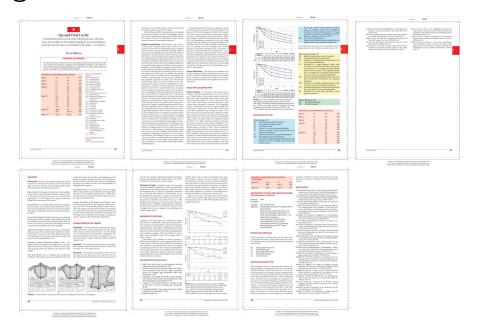
Preston, Wei, et al. "Towards Structuring Real-World Data at Scale: Deep Learning for Extracting Key Oncology Information from Clinical Text with Patient-Level Supervision", *Patterns 2023*.

## GPT-4: Structure Real-World Data

Preliminary results promising

"Read" annotation guideline → zero-shot structuring

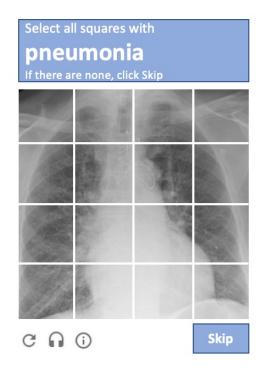




Project I	Hanover	Patient-Centri	ic Trial-Centric			Clinical Trial Triagi	ng					Cliff Wong
Accessio	on <b>No.:</b> 34 eb. 18, 19	DM JEFFREY I-234-58823 50		Sear	rch Builder Show	10 ¢ entries					Search:	
Gender: Histolog	y:	nocarcinoma)		₩	NCT No. ↑↓	Title ↑↓	Phase ↑↓	Matching Trial Diseases ↑↓	Matching Trial Stage ↑	Matching Trial Biomarkers	Notes ↑↓	Providence States
Stage Gr HLA type • HL/	oup: Stag e: A-A*02:01	HLA-A*02:01		⊗	NCT03953235	A Study of a Personalized Cancer Vaccine Targeting Shared Neoantigens	Phase 1/Phase 2	- Non-Small Cell Lung Carcinoma - Malignant Solid Neoplasm	- Metastatic - Advanced	- KRAS G12V	test3	CA, TX
		2 HLA-B*39:06 4 HLA-C*08:02		0	NCT04620330	A Study of Avutometinib (VS-6766) + Defactinib in Recurrent KRAS G12V, Other KRAS and BRAF Non-Small Cell Lung Cancer	Phase 2	- Non-Small Cell Lung Carcinoma		- KRAS G12V - KRAS Mutation	test6	CA, OR, TX
		ed Curation N/A		0	NCT03454035	Ulixertinib/Palbociclib in Patients With Advanced Pancreatic and Other Solid Tumors	Phase 1	- Malignant Solid Neoplasm	- Stage IV - Metastatic - Advanced	- KRAS G12X - KRAS Mutation		
Search Trial Fil	Repor	t		0	NCT05631899	Combination of CAR-DC Vaccine and Anti-PD-1 Antibody in Local Advanced/Metastatic Solid Tumors	Phase 1	- Malignant Solid Neoplasm	- Metastatic - Advanced	- KRAS G12V - KRAS Mutation		
☑ Age M		<u>.</u>		0	NCT05438667	TCR-T Cell Therapy on Advanced Pancreatic Cancer and Other Solid Tumors	Early Phase 1	- Malignant Solid Neoplasm	- Metastatic - Advanced	- KRAS G12V - KRAS Mutation		
	America			0	NCT04625647	Testing the Use of Targeted Treatment (AMG 510) for KRAS G12C Mutated Advanced Non-squamous Non-small Cell Lung Cancer (A Lung-MAP Treatment Trial)	Phase 2	- Non-Squamous Non- Small Cell Lung Carcinoma - Lung Adenocarcinoma - Non-Small Cell Lung Carcinoma - Lung Carcinoma	- Stage IVA - Stage IVB - Stage IV - Advanced	- KRAS Mutation		AK, CA, MT, NM, OR, TX, WA
□ Provid	ence State	es		⊗	NCT04999761	AB122 Platform Study	Phase 1	- Non-Squamous Non- Small Cell Lung Carcinoma - Non-Small Cell Lung	- Metastatic - Advanced	- KRAS Mutation		
clinical signif.	gene	protein change	variant					Carcinoma - Malignant Solid Neoplasm				
YES	KRAS	p.Gly12Val	G12V	0	NCT03667716	COM701 (an Inhibitor of PVRIG) in Subjects With Advanced Solid Tumors.	Phase 1	- Non-Small Cell Lung Carcinoma	- Stage IV - Metastatic - Advanced	- KRAS Mutation		CA, TX
YES	TP53	p.Arg306Ter	R306*					- Lung Carcinoma - Malignant Solid Neoplasm	. 10.1000			
YES	APC	p.Glu1353Ter p.Glu2139llefsTer6	E1353* E2139lfs*6	≪	NCT04511845	A Dose-Escalation Study of SPYK04 in Patients With Locally Advanced or Metastatic Solid Tumors (With Expansion).	Phase 1	- Non-Small Cell Lung Carcinoma	- Metastatic	- KRAS Mutation		TX
YES	ERBB2	3.4(fold-change)	ERBB2- High					- Malignant Solid Neoplasm		- MAPK/ERK pathway		16

## Biomedical LLMs

## General vs Health Labeled Data





### **IMPRESSION**

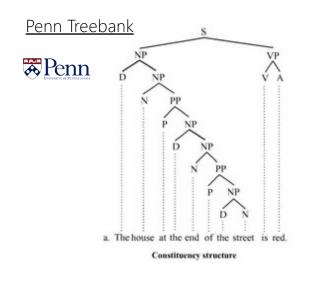
No significant change in right middle and low lobe pneumonia. Small increase in left pleural effusion. .....

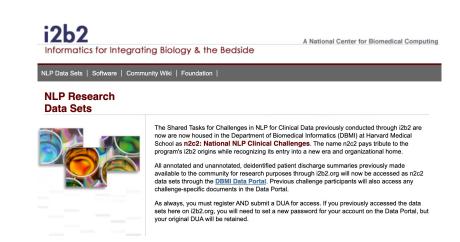


Two cows are grazing in the field.

Biomedical and clinical domain label require expertise

# General vs Health Data Availability

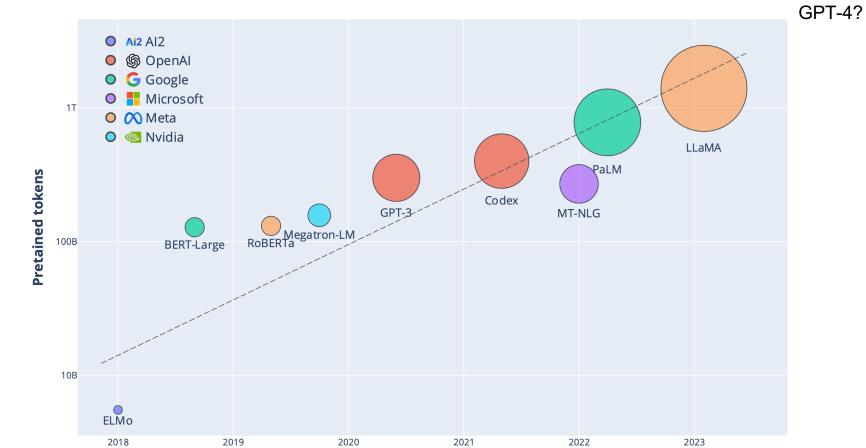




1992 2006

Comparable datasets over a decade later

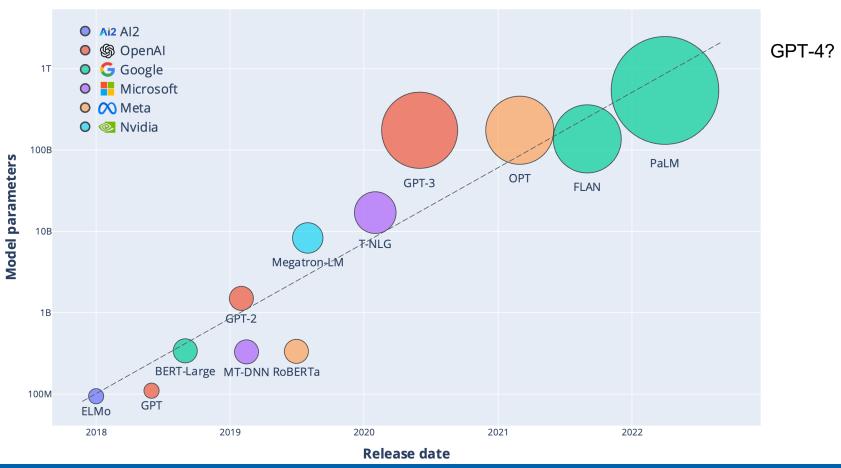
## Growth of Data (5B → 1T)

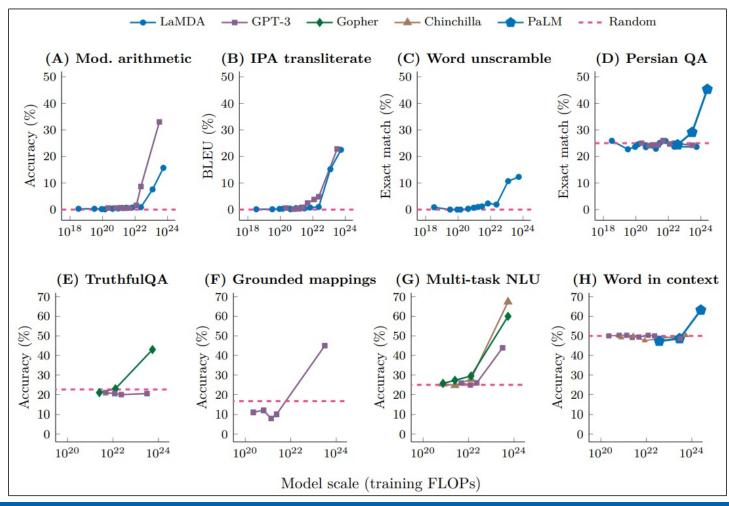


Microsoft Health Futures 20

**Release date** 

## Growth of Model Size (100M → 1T+)





Wei, et al. "Emergent Abilities of Large Language Models", *TMLR 2022*.

## Effects of Scale

350M

750M

3B

20B









A portrait photo of a kangaroo wearing an orange hoodie and blue sunglasses standing on the grass in front of the Sydney Opera House holding a sign on the chest that says Welcome Friends!

https://parti.research.google/

# General-purpose Interface

A1: visit summary

A2: PHI

A3: document type

A4: disease

. . . . .

Output interface

## LLM for language understanding and generation

Input interface

T1: summarize the doctor-patient dialogue.

T2: extract PHI from the patient note.

T3: classify the PubMed abstract.

T3: what disease does the patient have?

. . . . . .





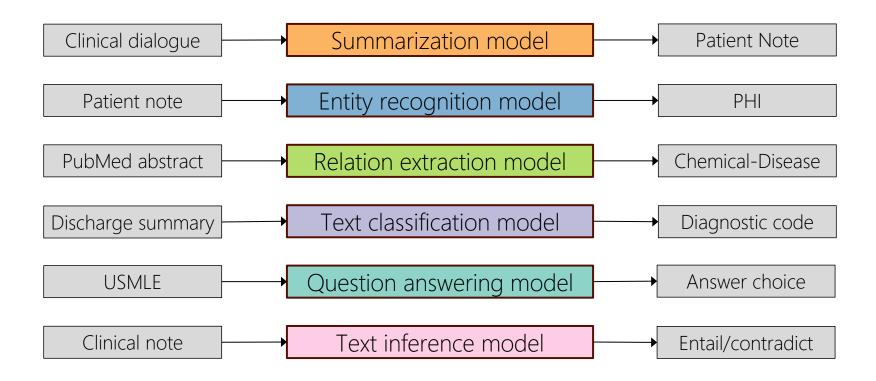




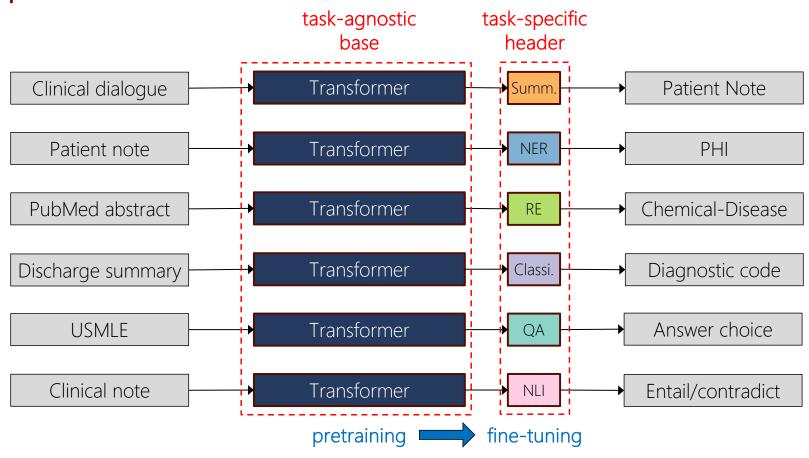
Closed-set Open-ended
Classification Generation

Representation Promptable Interface

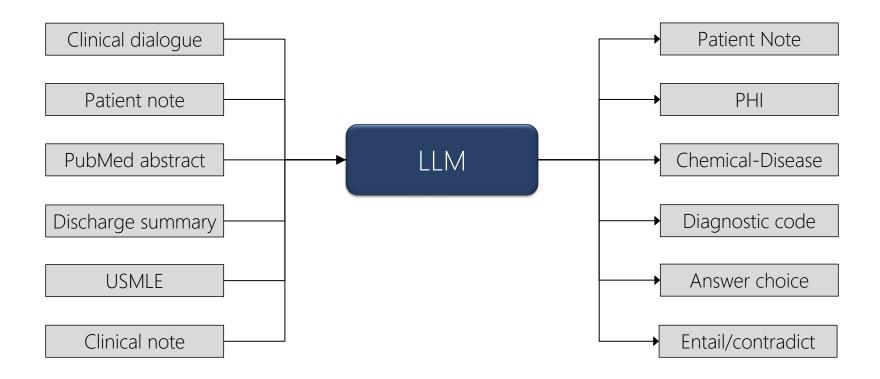
# **Specialist Models**



## **Specialist Headers**



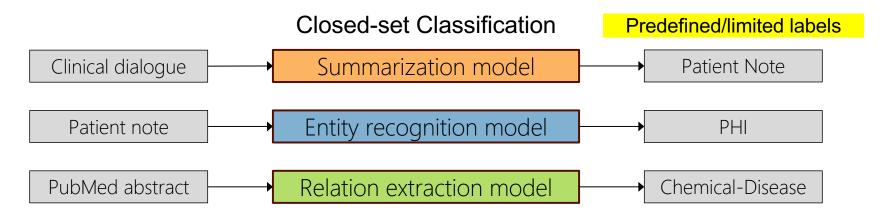
## **Generalist Models**

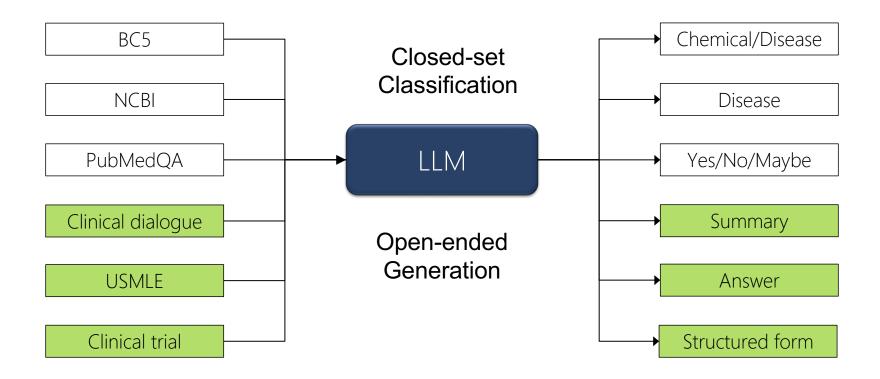






Representation Promptable Interface





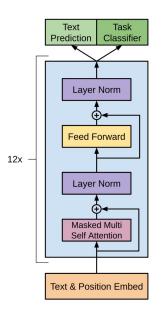






### Representation learning

- Expensive
- Engineering heavy
- Task-specific



### Promptable interface

- Training free
- Universal interface natural language



Improving Language Understanding by Generative Pre-Training
Retrieval-based Language Models and Applications

## **Biomedical LLMs**





















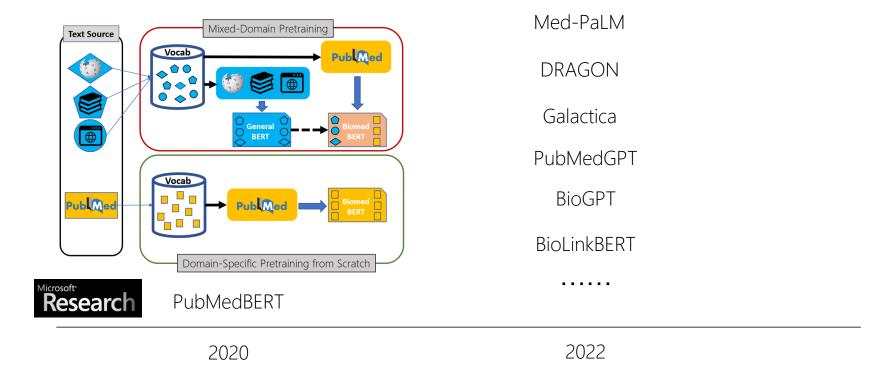




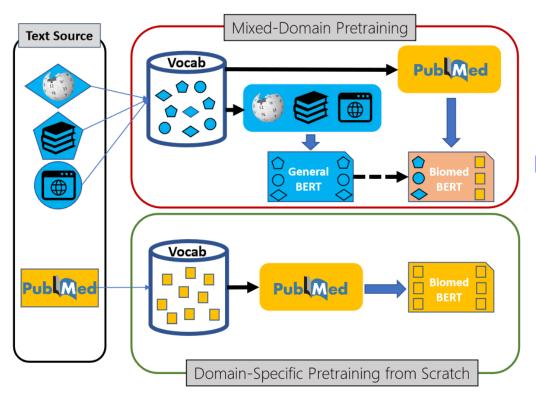




# Domain-Specific Pretraining



# Why Domain-Specific Pretraining?



Yu, et al. "Domain-Specific Language Model Pretraining for Biomedical Natural Language Processing", Special Issue on Computational Methods for Biomedical Natural Language Processing, ACM Transactions on Computing for Health 2021.

### BioMedBERT (formerly PubMedBERT)

In **bounded-resource** scenarios, enable **more efficient learning** by focusing on in-domain data

# Why Domain-Specific Pretraining?

### Shattered into pieces

# Domain-specific Vocab

Biomedical Term	Category	BERT	SciBERT	PubMedBERT (Ours)
diabetes	disease	✓	✓	<u> </u>
leukemia	disease	✓	$\checkmark$	✓
lithium	drug	$\checkmark$	$\checkmark$	✓
insulin	drug	✓	✓	✓
DNA	gene	$\checkmark$	$\checkmark$	✓
promoter	gene	✓	✓	$\checkmark$
hypertension	disease	hyper-tension	$\checkmark$	✓
nephropathy	disease	ne-ph-rop-athy	$\checkmark$	$\checkmark$
lymphoma	disease	l-ym-ph-oma	$\checkmark$	<u>I</u>
lidocaine	drug	lid-oca-ine]	$\checkmark$	$\checkmark$
oropharyngeal	organ	oro-pha-ryn-ge-al	or-opharyngeal	· ✓
cardiomyocyte	cell	card-iom-yo-cy-te	cardiomy-ocyte	$\checkmark$
chloramphenicol	drug	ch-lor-amp-hen-ico-l	chlor-amp-hen-icol	✓
RecA	gene	Rec-A	Rec-A	✓
acetyltransferase	gene	ace-ty-lt-ran-sf-eras-e	acetyl-transferase	✓
clonidine	drug	cl-oni-dine	clon-idine	✓
naloxone	drug	na-lo-xon-e	nal-oxo-ne	✓

Yu, et al. "Domain-Specific Language Model Pretraining for Biomedical Natural Language Processing", Special Issue on Computational Methods for Biomedical Natural Language Processing, ACM Transactions on Computing for Health 2021.

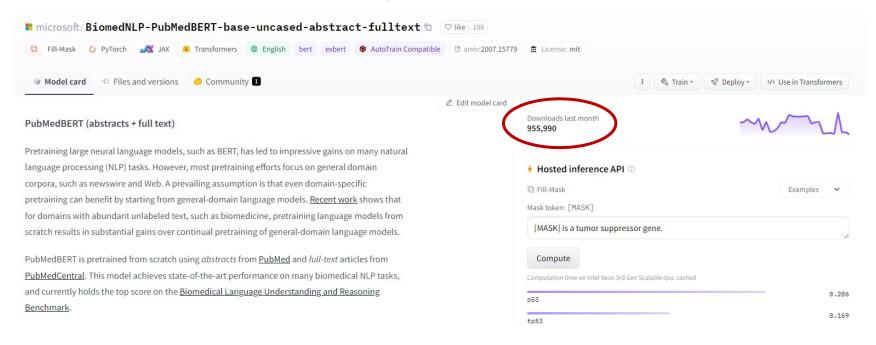
### Domain-specific Vocab

### Preserves the integrity of

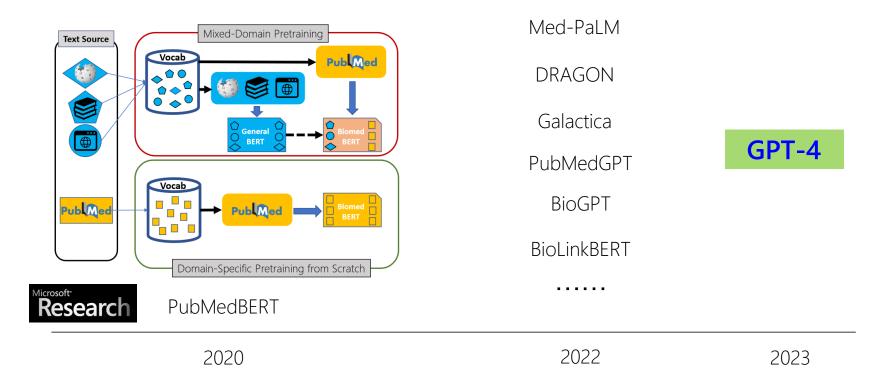
- Biomedical terms
- Amino acid sequences
- SMILES formula
- DNA sequences
- Mathematics
- Citations

• etc.

# BioMedBERT (formerly PubMedBERT)



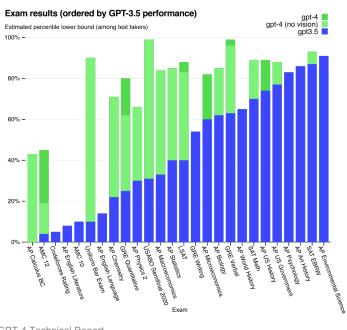
# Domain-Specific Pretraining → Generalist Model



## GPT-4

## Out-of-Box: Expert-Level Competency on USMLE

The most powerful general-purpose LLM Human-level performance on many tasks



- SOTA on medical competency examinations
- "How well does the AI perform clinically? And my answer is, I'm stunned to say: Better than many doctors I've observed." — Isaac Kohane MD

D	GPT-4-base	GPT-4
Dataset	$5~\mathrm{shot}$ / $0~\mathrm{shot}$	$5~\mathrm{shot}$ / $0~\mathrm{shot}$
MedQA		
Mainland China	<b>78.63</b> / 74.34	75.31 / 71.07
Taiwan	<b>87.47</b> / 85.14	84.57 / 82.17
US (5-option)	<b>82.25</b> / 81.38	78.63 / 74.71
US (4-option)	<b>86.10</b> / 84.45	81.38 / 78.87
PubMedQA		
Reasoning Required	77.40 / 80.40	$74.40 \ / \ 75.20$
MedMCQA		
Dev	<b>73.66</b> / 73.42	$72.36 \ / \ 69.52$



**GPT-4 Technical Report** 

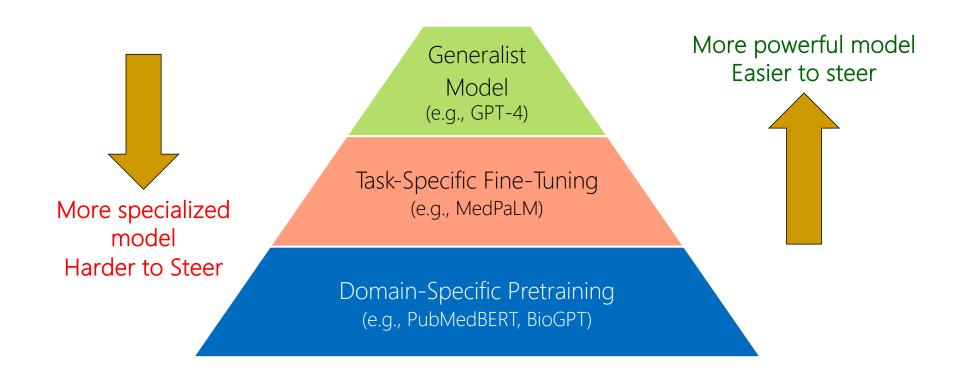
Capabilities of GPT-4 on Medical Challenge Problems The Al Revolution in Medicine: GPT-4 and Beyond

## GPT-4

GPT-4 has been pretrained on a large portion of the public web, which already contains a lot of biomedical text.

Component	Raw Size
Pile-CC	227.12 GiB
PubMed Central	90.27 GiB
Books3 <sup>†</sup>	100.96 GiB
OpenWebText2	62.77 GiB
ArXiv	56.21 GiB
Github	95.16 GiB
FreeLaw	51.15 GiB
Stack Exchange	32.20 GiB
<b>USPTO Backgrounds</b>	22.90 GiB
PubMed Abstracts	19.26 GiB
Gutenberg (PG-19) <sup>†</sup>	10.88 GiB
OpenSubtitles <sup>†</sup>	12.98 GiB
Wikipedia (en) <sup>†</sup>	6.38 GiB
DM Mathematics <sup>†</sup>	7.75 GiB
Ubuntu IRC	5.52 GiB
BookCorpus2	6.30 GiB
EuroParl <sup>†</sup>	4.59 GiB
HackerNews	3.90 GiB
YoutubeSubtitles	3.73 GiB
PhilPapers	2.38 GiB
NIH ExPorter	1.89 GiB
Enron Emails†	0.88 GiB
The Pile	825.18 GiB

# Generalist Models: Superior Steerability



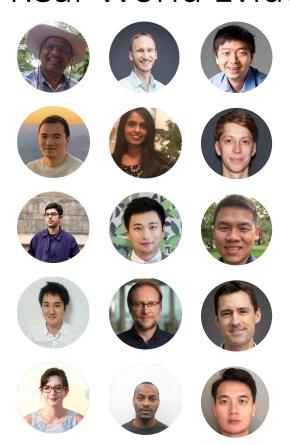
## Population-Level Health LLM

Patient → Serialized multimodal token sequence Initialize: GPT-101 (consumed entire public web) Continued pretraining: 8 billion "health documents"

What is the multimodal health scaling law?

Will there be emergent health capabilities?

## Real-World Evidence



- JAX: Susan Mockus, Sara Patterson
- Fred Hutchinson: Christopher Li, Kathi Malone
- Providence: Carlo Bifulco, Brian Piening
- MSR: Xiaodong Liu, Hao Cheng, Jianfeng Gao, Ozan Oktay, Javier Alvarez-Valle, Naveen Valluri
- Interns: Maxim Grechkin, Ankur Parikh, Victoria Lin, Sheng Wang, Stephen Mayhew, Daniel Fried, Violet Peng, Hai Wang, Robin Jia, Matthew McDermott, Alexis Ross, Zelalem Gero, Sarthak Jain, Jenny Chen, Hunter Lang, Benedikt Boecking, Varsha Kishore, Jinfeng Xiao, Wenxuan Zhou, Neha Hulkund, Michelle Li, Peniel Argaw, Hanwen Xu, Mars Huang, Juan Manuel Zambrano Chaves, Yiqing Xie, Isabel Chien, Alicia Curth