

DiReCT: Diagnostic Reasoning for Clinical Notes via Large Language Models

Bowen Wang, Jiuyang Chang, Yiming Qian, Guoxin Chen, Junhao Chen, Zhouqiang Jiang,
Jiahao Zhang, Yuta Nakashima, Hajime Nagahara

Osaka University

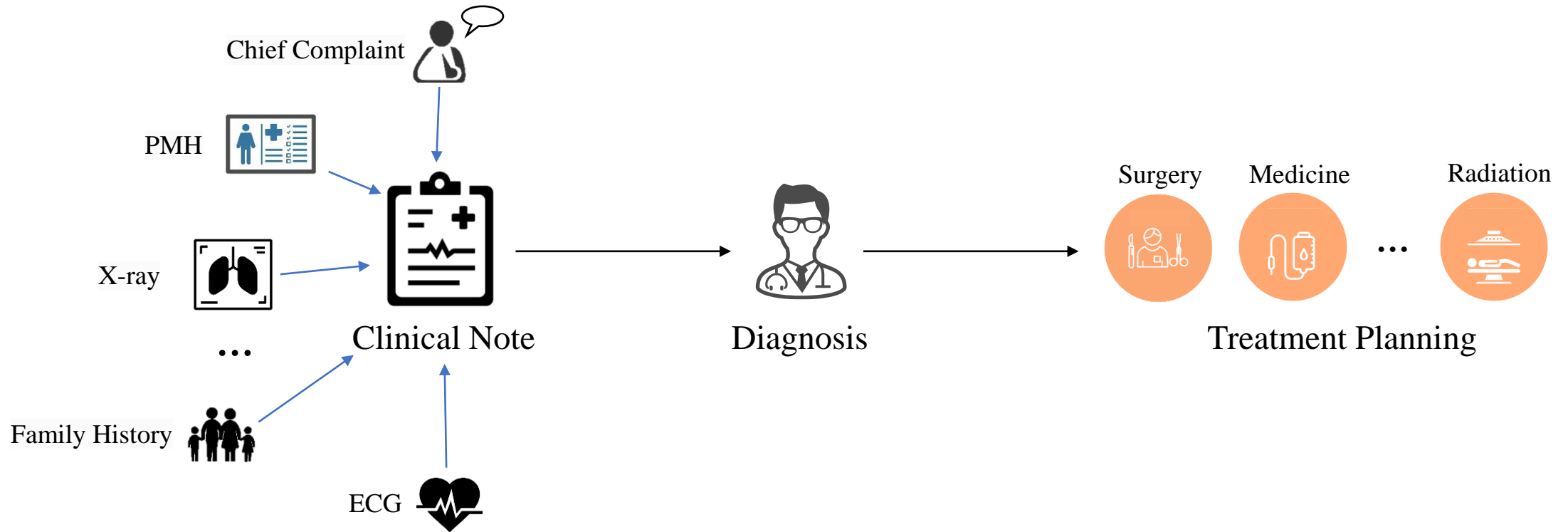
The First Affiliated Hospital of Dalian Medical University

Agency for Science, Technology and Research

Institute of Computing Technology, Chinese Academy of Science

Overview

Correct analysis of clinical records is crucial for treatment planning.



Overview

Complex and lengthy records can overwhelm doctors and increase diagnostic errors.

Chief Complaint: Scrotal and leg swelling ...
History of Present Illness: ... In the last 3 days his ***** has become quite swollen. It is similar ***** swelling when he was admitted with acute CHF ... EKG was consistent with priors (NSR, NANI, ***** changes). The left ventricle is mildly enlarged. He was given ***** with good UOP ...
Past Medical History: ... -Diabetes, -Hypertension, -CKD, stage 3, -GERD, -Depression, - Amputation of ***** , Pneumonia, - Osteoarthritis- History of ***** , Asthma ...
Family History: There is no family history of ***** artery ...
Physical Exam: ... LUNG: bibasilar rales that do not clear with deep inspiration. ... ABDOMEN: nondistended, ***** all quadrants. EXTREMITIES: bilateral pitting edema to the sacrum, extending to the up abdomen. Warm, well perfused. ... HEENT: AT/NC, EOMI, PERRL. ...
Pertinent Results: __ 03:50PM BLOOD WBC-8.0 RBC-3.26* Hgb-9.3* Hct-30.9* MCHC-29.9* ... __ 11:30AM BLOOD proBNP-3843 ... Overall left ventricular systolic function is mildly depressed (LVEF= 45-50 %) without regional wall motion abnormalities. ***** imaging suggests an increased ***** filling pressure (PCWP>*****Hg) ...

Clinical Note



Records are complex
Lack of specialist



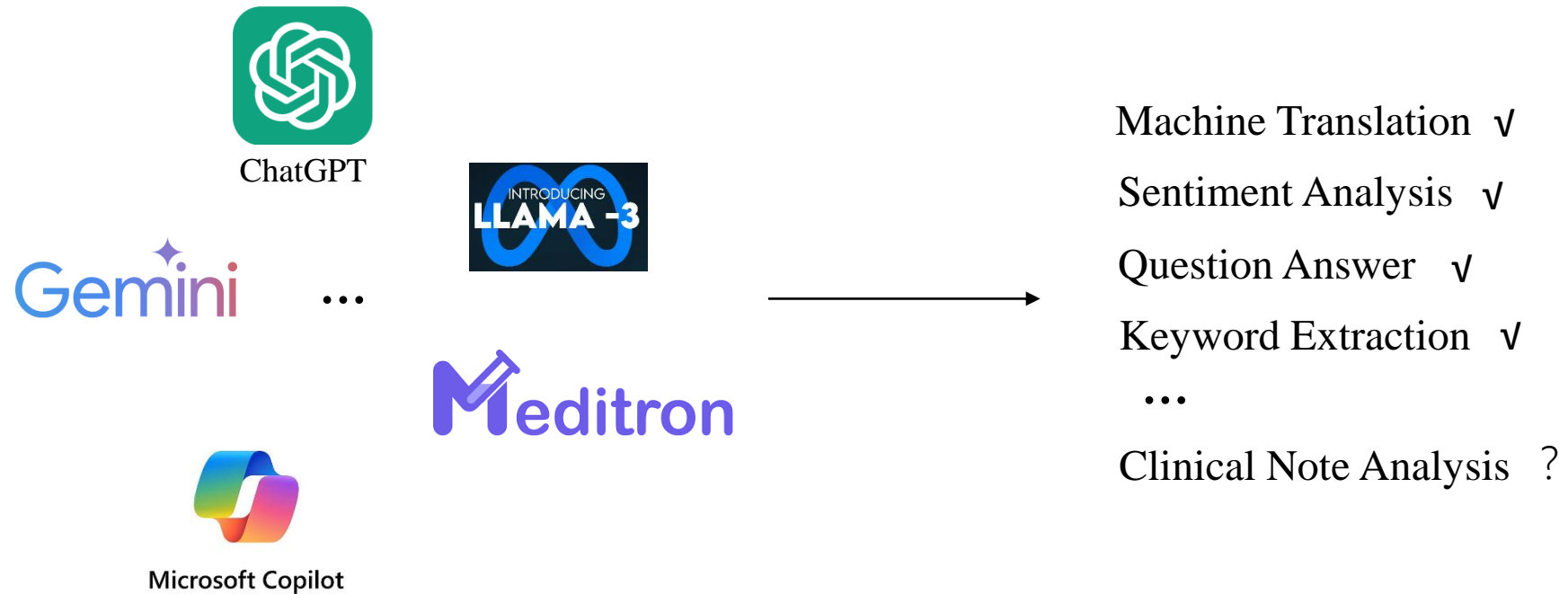
Misdiagnosis



Automatic diagnosis is necessary

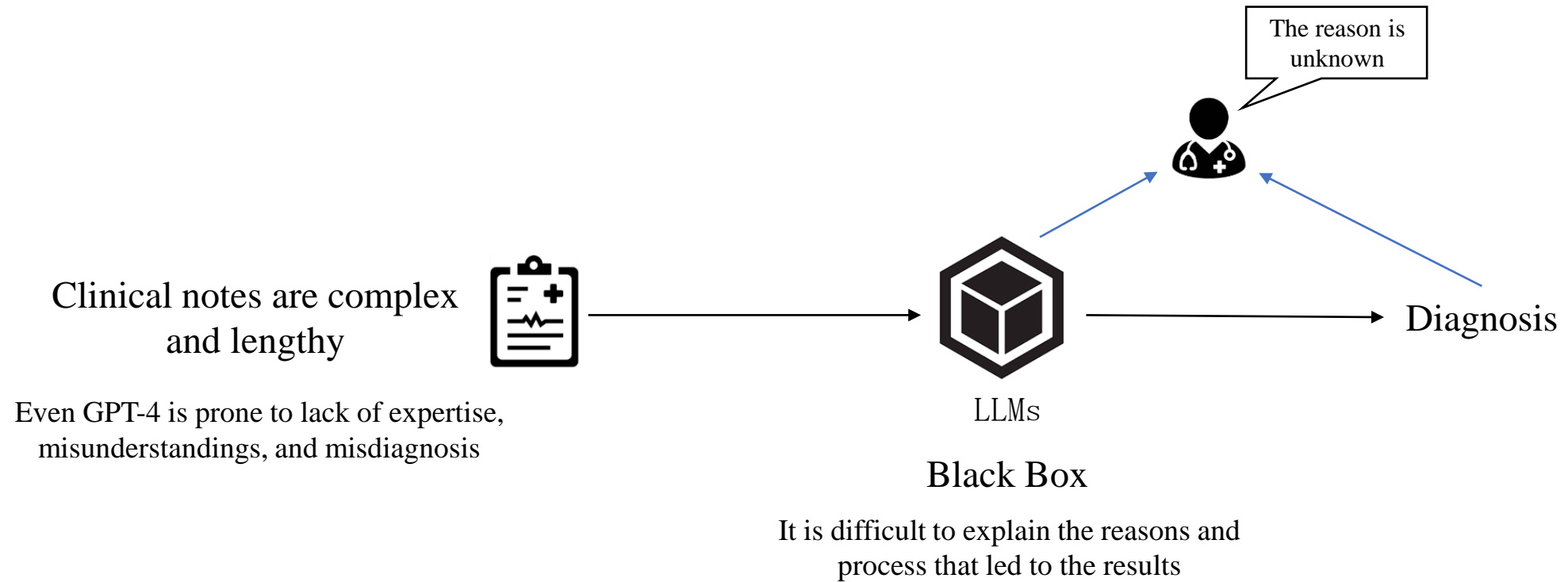
Overview

Recently, large language models (LLMs) have shown their power in a variety of language tasks.



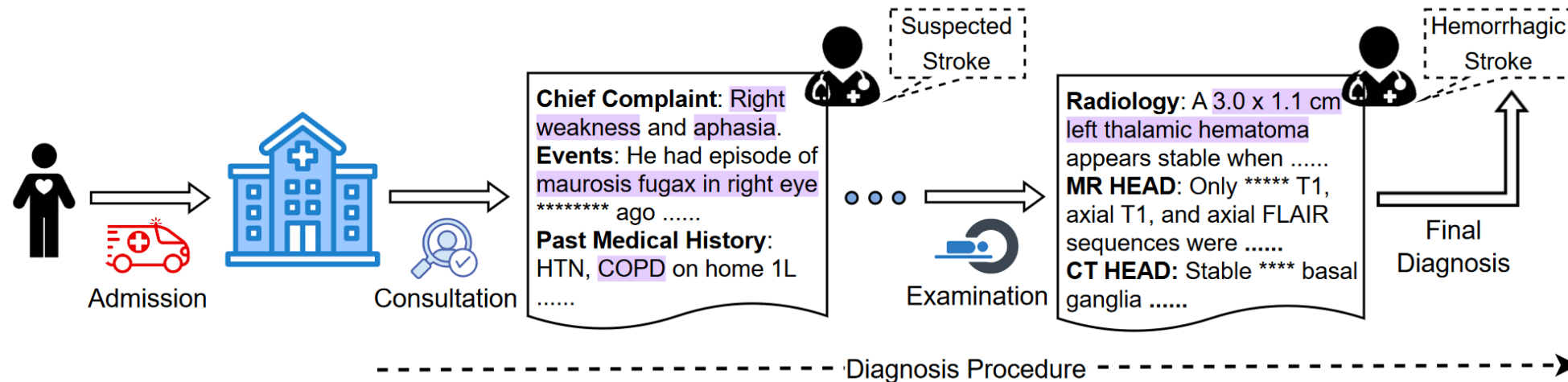
Overview

It is important that LLMs are explainable and consistent with physicians.



Overview

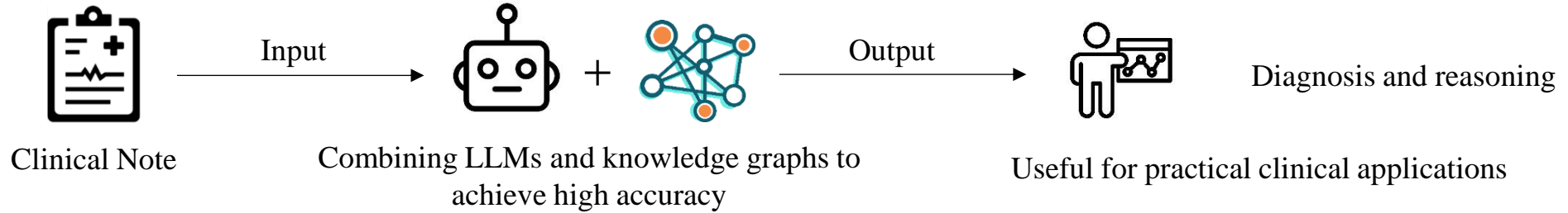
The diagnostic process of a human doctor follows existing diagnostic rules.



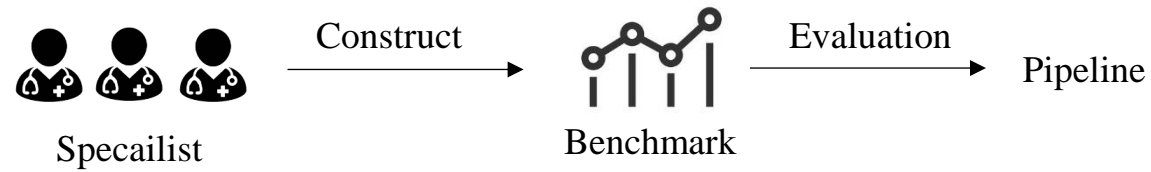
When a patient is admitted, an initial consultation takes place to collect subjective information. Subsequent observations may then require further examination to confirm the diagnosis.

Overview

An interpretable pipeline

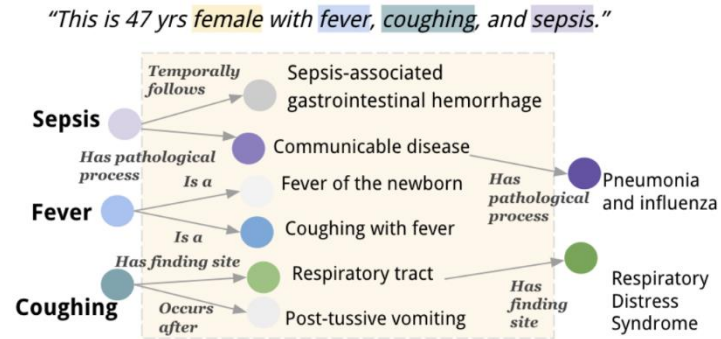


A benchmark dataset



Diagnostic Knowledge Graph

UMLS Knowledge Graph (existing)



Our Diagnostic Knowledge Graph

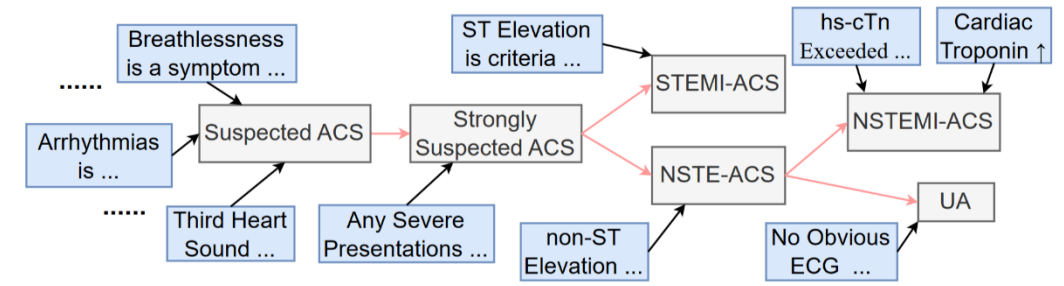


Figure 2: A part of k_i for i being Acute Coronary Syndromes.

- Only provide relations for simple words (For use this KG, input has to be split into words.)
- No experimental values.

$$g_i = (\mathcal{D}_i, \mathcal{F}_i)$$

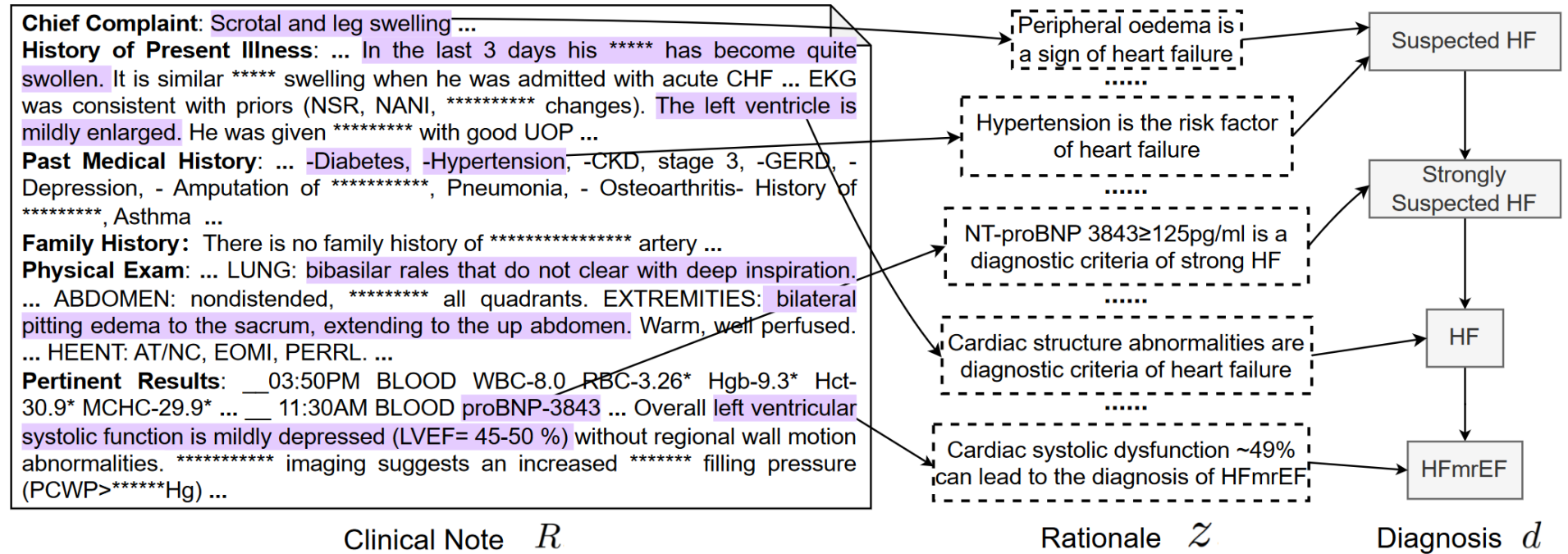
$$k_i = (\mathcal{D}_i, \mathcal{P}_i, \mathcal{S}_i, \mathcal{F}_i)$$

\mathcal{D}_i Diagnostic Nodes \mathcal{F}_i Procedural Edges

\mathcal{P}_i Premise Nodes \mathcal{S}_i Supporting Edges

\mathcal{D}^* Leaf Nodes

Data Annotation



An annotation sample of Heart Failure (HF). The left part is the clinical note alongside extracted observations by a doctor. The middle part outlines the steps of the rationale for the premise corresponding to each diagnostic node shown in the right part.

$$\mathcal{E} = \{(o, z, d)\} \text{ for } (R, d^*)$$

Annotated by 9 clinical physicians and subsequently verified for accuracy and completeness by three senior medical experts.

Data Statistics and Task Definition

Statistics of all 25 disease categories, 511 annotations.

Domains	Categories	# samples	$ \mathcal{D}_i $	$ \mathcal{D}_i^* $
Cardiology	Acute Coronary Syndromes	65	6	3
	Aortic Dissection	14	3	2
	Atrial Fibrillation	10	3	2
	Cardiomyopathy	9	5	4
	Heart Failure	52	6	3
	Hyperlipidemia	2	2	1
	Hypertension	32	2	1
Gastroenterology	Gastritis	27	5	3
	Gastroesophageal Reflux Disease	41	2	1
	Peptic Ulcer Disease	28	3	2
	Upper Gastrointestinal Bleeding	7	2	1
Neurology	Alzheimer	10	2	1
	Epilepsy	8	3	2
	Migraine	4	3	2
	Multiple Sclerosis	27	6	4
	Stroke	28	3	2
Pulmonology	Asthma	33	7	5
	COPD	19	6	4
	Pneumonia	20	4	2
	Pulmonary Embolism	15	5	3
	Tuberculosis	5	3	2
Endocrinology	Adrenal Insufficiency	20	4	3
	Diabetes	13	4	2
	Pituitary	12	4	3
	Thyroid Disease	10	6	4

Statistics of 5 medical domains.

Medical domain	# cat.	# samples	$ \mathcal{D}_i $	$ \mathcal{D}_i^* $	$ \mathcal{O} $	Length
Cardiology	7	184	27	16	8.7	1156.6 t
Gastroenterology	4	103	11	7	4.3	1026.0 t
Neurology	5	77	17	11	11.9	1186.3 t
Pulmonology	5	92	26	17	10.7	940.7 t
Endocrinology	4	55	20	14	6.9	1063.5 t
Overall	25	521	101	65	8.5	1074.6 t

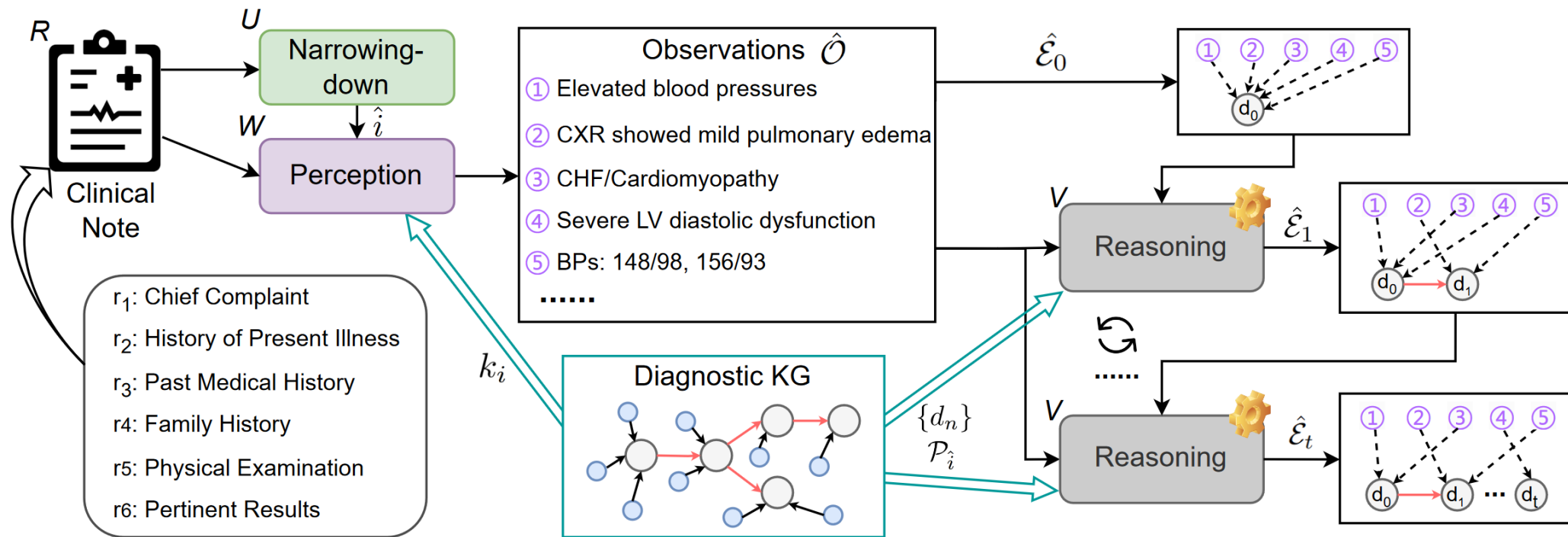
Each note only has a primary discharge diagnosis (PDD) in \mathcal{D}^*

$$\text{Task 1 } \hat{d}^*, \hat{\mathcal{E}} = M(R, \mathcal{G})$$

$$\text{Task 2 } \hat{d}^*, \hat{\mathcal{E}} = M(R, \mathcal{K})$$

An AI Agent Pipeline

Our baseline comprises three LLM-based modules: narrowing-down U , perception W , and reasoning V .



Experiment Results

- Accuracy of diagnosis (Acc)
- Completeness of observations (Obs)
- Faithfulness of explanations (Exp)

Auto Evaluation via LLama3 8B

Table 3: Diagnostic reasoning ability of different LLMs under the proposed baseline method.

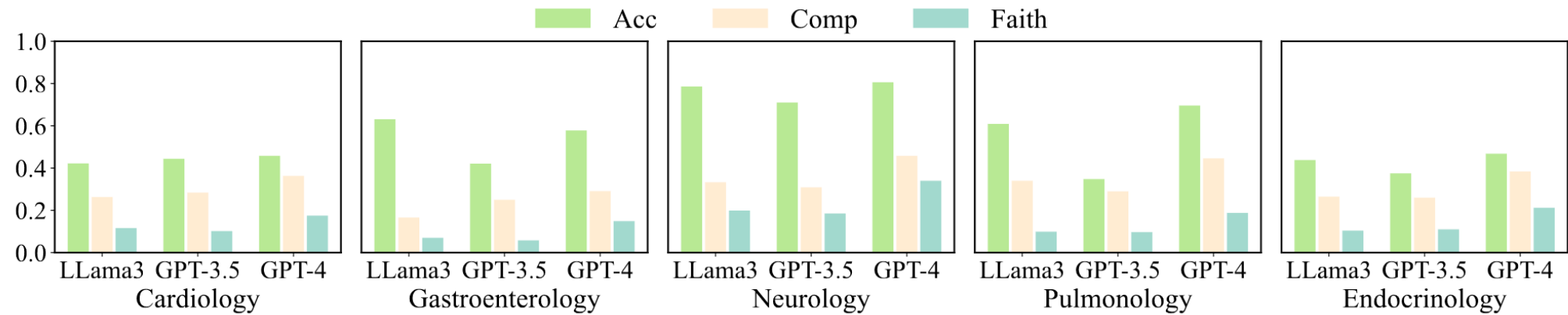
Task	Models	Diagnosis		Observation			Explanation	
		Acc^{cat}	Acc^{diag}	Obs^{pre}	Obs^{rec}	Obs^{comp}	Exp^{com}	Exp^{all}
With \mathcal{G}	Zephyr 7B	0.274	0.151	0.123 \pm 0.200	0.115 \pm 0.166	0.092 \pm 0.108	0.071 \pm 0.139	0.014 \pm 0.037
	Mistral 7B	0.507	0.306	0.211 \pm 0.190	0.317 \pm 0.253	0.173 \pm 0.157	0.230 \pm 0.312	0.062 \pm 0.088
	Mixtral 8 \times 7B	0.413	0.237	0.147 \pm 0.165	0.266 \pm 0.261	0.124 \pm 0.138	0.144 \pm 0.268	0.029 \pm 0.056
	LLama3 8B	0.576	0.321	0.253 \pm 0.156	0.437 \pm 0.207	0.219 \pm 0.137	0.232 \pm 0.316	0.071 \pm 0.093
	LLama3 70B	0.752	0.540	0.277 \pm 0.146	0.537 \pm 0.192	0.256 \pm 0.142	0.395 \pm 0.320	0.112 \pm 0.110
	GPT-3.5 turbo	0.679	0.455	0.389 \pm 0.212	0.351 \pm 0.192	0.275 \pm 0.167	0.331 \pm 0.366	0.103 \pm 0.127
	GPT-4 turbo	0.772	0.572	0.446 \pm 0.207	0.491 \pm 0.180	0.371 \pm 0.186	0.475 \pm 0.363	0.199 \pm 0.181
With \mathcal{K}	LLama3 8B	0.576	0.344	0.235 \pm 0.162	0.394 \pm 0.227	0.199 \pm 0.142	0.327 \pm 0.375	0.087 \pm 0.114
	LLama3 70B	0.735	0.581	0.262 \pm 0.146	0.501 \pm 0.208	0.236 \pm 0.131	0.463 \pm 0.374	0.125 \pm 0.117
	GPT-3.5 turbo	0.652	0.413	0.347 \pm 0.241	0.279 \pm 0.203	0.232 \pm 0.184	0.374 \pm 0.408	0.121 \pm 0.152
	GPT-4 turbo	0.781	0.614	0.431 \pm 0.207	0.458 \pm 0.187	0.353 \pm 0.170	0.633 \pm 0.338	0.247 \pm 0.201

Table 4: Evaluation of diagnostic reasoning ability of LLMs when no external knowledge is provided.

Task	Models	Acc^{diag}	Observation			Explanation	
			Obs^{pre}	Obs^{rec}	Obs^{comp}	Exp^{com}	Exp^{all}
With \mathcal{D}^*	LLama3 8B	0.070	0.154 \pm 0.139	0.330 \pm 0.244	0.135 \pm 0.122	0.020 \pm 0.100	0.004 \pm 0.016
	LLama3 70B	0.502	0.257 \pm 0.150	0.509 \pm 0.213	0.237 \pm 0.145	0.138 \pm 0.209	0.034 \pm 0.054
	GPT-3.5 turbo	0.223	0.164 \pm 0.242	0.149 \pm 0.212	0.116 \pm 0.174	0.091 \pm 0.231	0.025 \pm 0.065
	GPT-4 turbo	0.636	0.461 \pm 0.206	0.482 \pm 0.160	0.378 \pm 0.174	0.186 \pm 0.221	0.074 \pm 0.090
No Knowledge	LLama3 8B	0.023	0.137 \pm 0.159	0.258 \pm 0.274	0.119 \pm 0.141	0.018 \pm 0.083	0.006 \pm 0.026
	LLama3 70B	0.037	0.246 \pm 0.148	0.504 \pm 0.222	0.227 \pm 0.148	0.022 \pm 0.093	0.007 \pm 0.030
	GPT-3.5 turbo	0.059	0.161 \pm 0.238	0.148 \pm 0.215	0.113 \pm 0.171	0.036 \pm 0.131	0.011 \pm 0.039
	GPT-4 turbo	0.074	0.410 \pm 0.208	0.443 \pm 0.191	0.324 \pm 0.182	0.047 \pm 0.143	0.019 \pm 0.058

Experiment Results

Performance of LLama3 70B, GPT-3.5, and GPT-4 under different medical domains.



Experiment Results

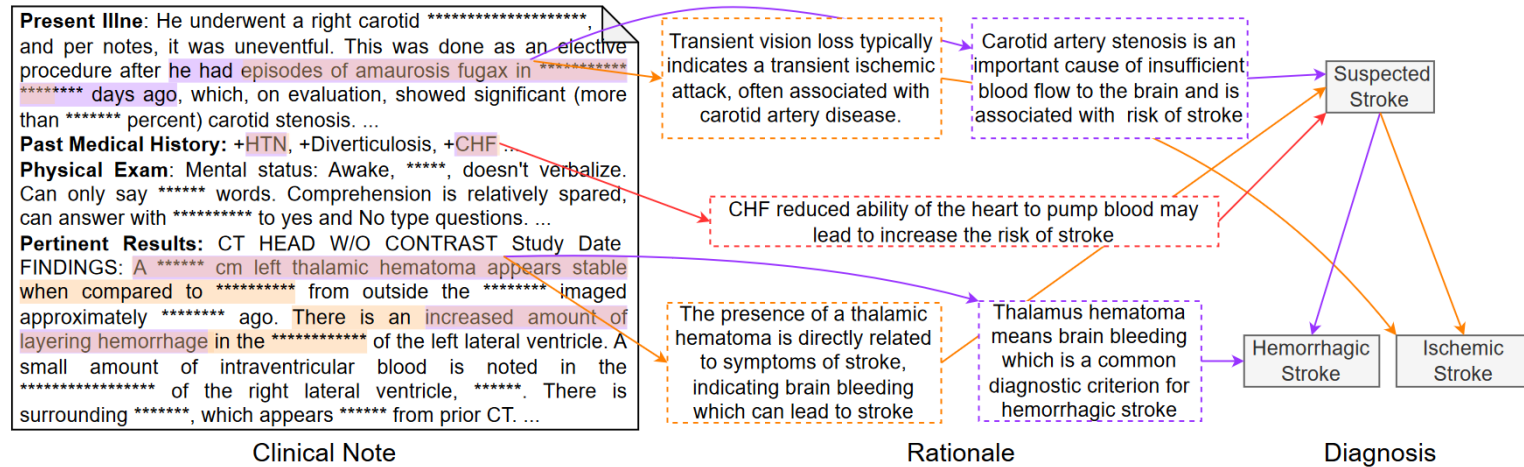
Consistency of automated evaluation metrics with human judgments.

Table 5: Consistency of automated evaluation with human judgments. Evaluated by mean and confidence interval (CI).

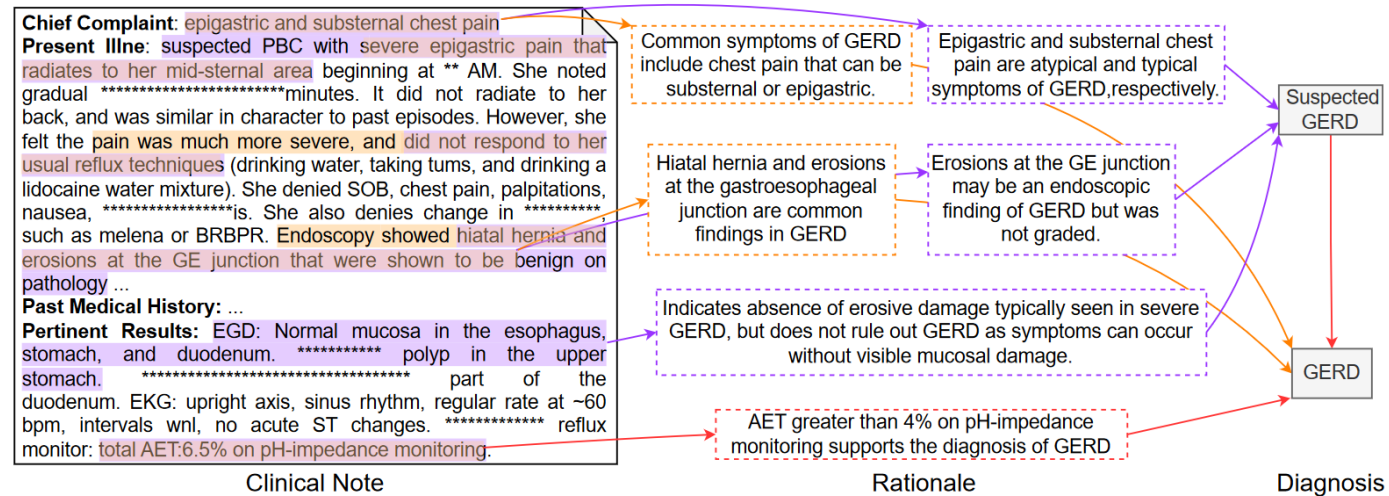
Model	Observation		Rationalization	
	Mean	95% CI	Mean	95% CI
LLama3 8B	0.887	0.844 ~ 0.878	0.835	0.759 ~ 0.818
GPT-4 turbo	0.902	0.830 ~ 0.863	0.876	0.798 ~ 0.853

Generated Samples From GPT-4

Purple, orange, and red indicate ground truth, prediction, and common in both, respectively.



An example prediction for a clinical note with PDD of Hemorrhagic Stroke by GPT-4.



An example prediction for a clinical note with PDD of GERD by GPT-4

Thanks