

Introduction

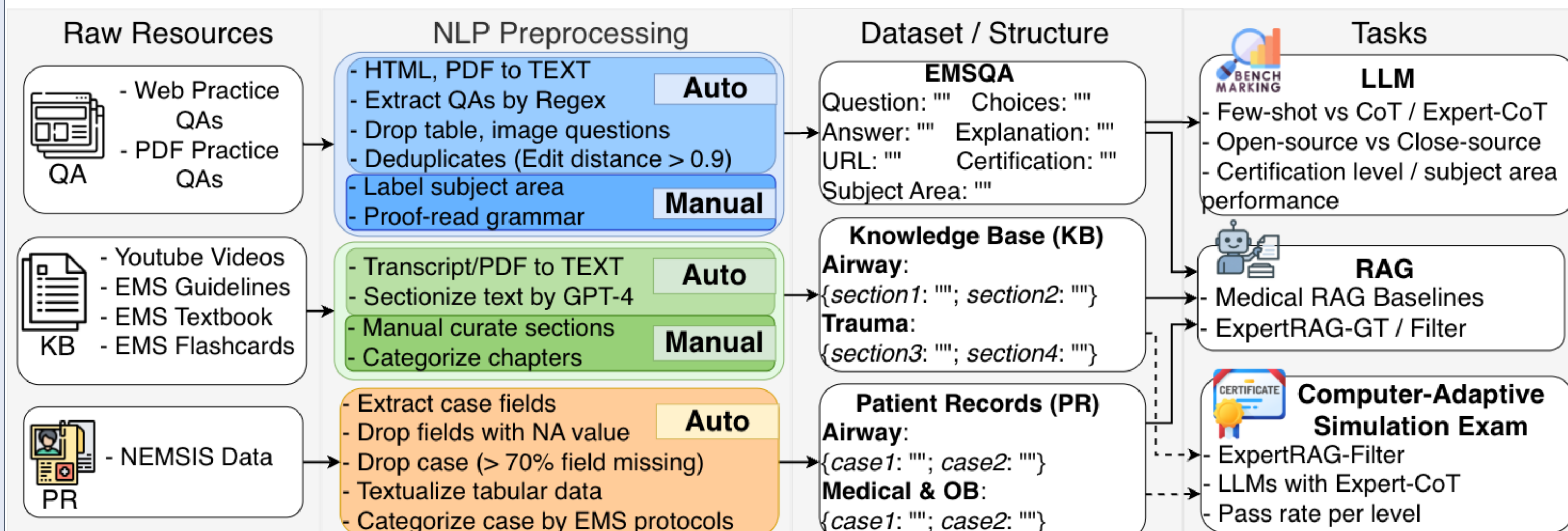
❖ Motivation

- **High-Stake Need:** Emergency Medical Service decisions require reliable, expert-level reasoning
- **Method Gap:** Current CoT/RAG treat reasoning/retrieval as generic, ignoring which expertise (subject area + certification) should guide the process.
- **Data Gap:** Existing medical QA lack structured expertise annotations and EMS-specific curated knowledge, making expert-aligned evaluation hard.

	MedQA	MedMCQA	EMSQA
Domain	General Med.	General Med.	EMS
Data Size	12.7K	193K	24.3K
Exam	USMLE	AIIMS&NEET PG	NREMT
#Certification	1	1	4
#Subject Area	X	21	10
KB	Raw	X	Categorized

❖ Contribution

- **EMSQA:** 24.3K EMS-related QAs, covering 10 subject areas and 4 certifications. Accompanied by **EMSKB** with 40K documents and 4M Real-world patient care reports
- **Expertise-guided LLM framework**, infer EMS expertise attributes and inject them via 1) **Expert-CoT** (expert-style reasoning prompts) and 2) **ExpertRAG** (expert-aligned retrieval from EMS KBs + patient records).
- Expert-CoT + ExpertRAG improves accuracy by up to **4.59%**; 32B expertise-augmented models pass all **EMS certification simulations**.



EMSQA, EMSKB and Patient Records

Set	Size	Type	Criteria	vs KB	vs PR
Public	Train(13,021)	Semantic	Avg Sim	79.21	66.45
	Val(1,860)	Semantic	Vocab	82.95	21.14
	Test(3,721)	Syntactic (hit rate)	Cpt w/o norm	41.65	8.87
		Syntactic (hit rate)	Cpt w/ norm	63.30	15.28
Private	Test(5,669)	Semantic	Avg Sim	80.75	75.35
		Semantic	Vocab	90.89	28.26
		Syntactic (hit rate)	Cpt w/o norm	53.18	14.36
		Syntactic (hit rate)	Cpt w/ norm	72.49	22.66

Data	Split	#Explanations	#Choices (avg/max)	#Answers (avg/max)	Question Tokens (avg/max)	Choice Tokens (avg/max)	Tokens	Vocab
Public	Train (13,021)	2217	4.01 / 7.00	1.00 / 3.00	18.27 / 218	6.28 / 240	565,303	14,017
	Val (1,860)	383	3.99 / 5.00	1.00 / 3.00	19.12 / 155	6.01 / 44	80,215	6,629
	Test (3,721)	773	4.01 / 6.00	1.00 / 3.00	18.99 / 135	6.10 / 60	161,464	8,913
	Total (18,602)	3132	4.01 / 7.00	1.00 / 3.00	18.50 / 218	6.22 / 240	806,982	16,032
Private	Test (5,669)	5451	4.01 / 6.00	1.06 / 4.00	30.44 / 355	5.46 / 47	296,673	10,637

❖ EMSQA Statistics

- 18,602 public and 5,669 private NREMT practice questions
- 4 Certification levels, 10 Subject Areas

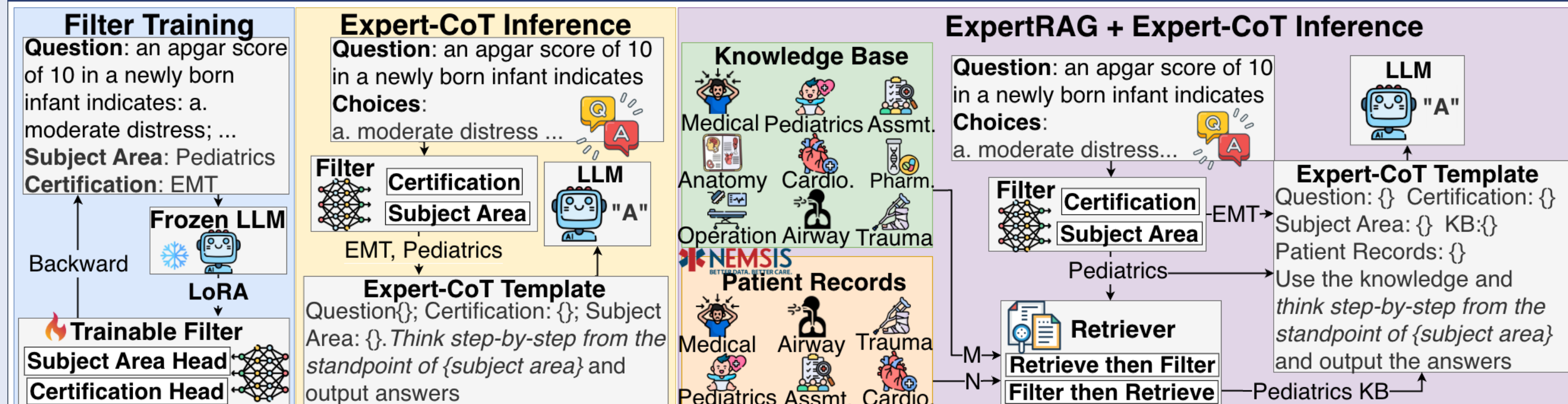
❖ EMSKB

- Crawled from 16 education resources
- Knowledge Coverage
 - Syntactic (Embedding Similarity)
 - Semantic (Vocab, EMS Concept)

❖ Patient Records

- NEMSIS Tubular Data → Textual Data
- Knowledge Coverage (Syntactic, Semantic)

Expert-CoT & ExpertRAG



❖ Filter Training

- **Input:** Q + <classify> **Output:** Subject Area; Certification

❖ Expert-CoT

$$\hat{A}_i = f^{\text{CoT-Expert}}(q_i, O_i, \hat{l}_i, \hat{s}_i).$$

- Given the question q_i , options O_i , certification level \hat{l}_i , Think from EMS-specific standpoint of subject area \hat{s}_i

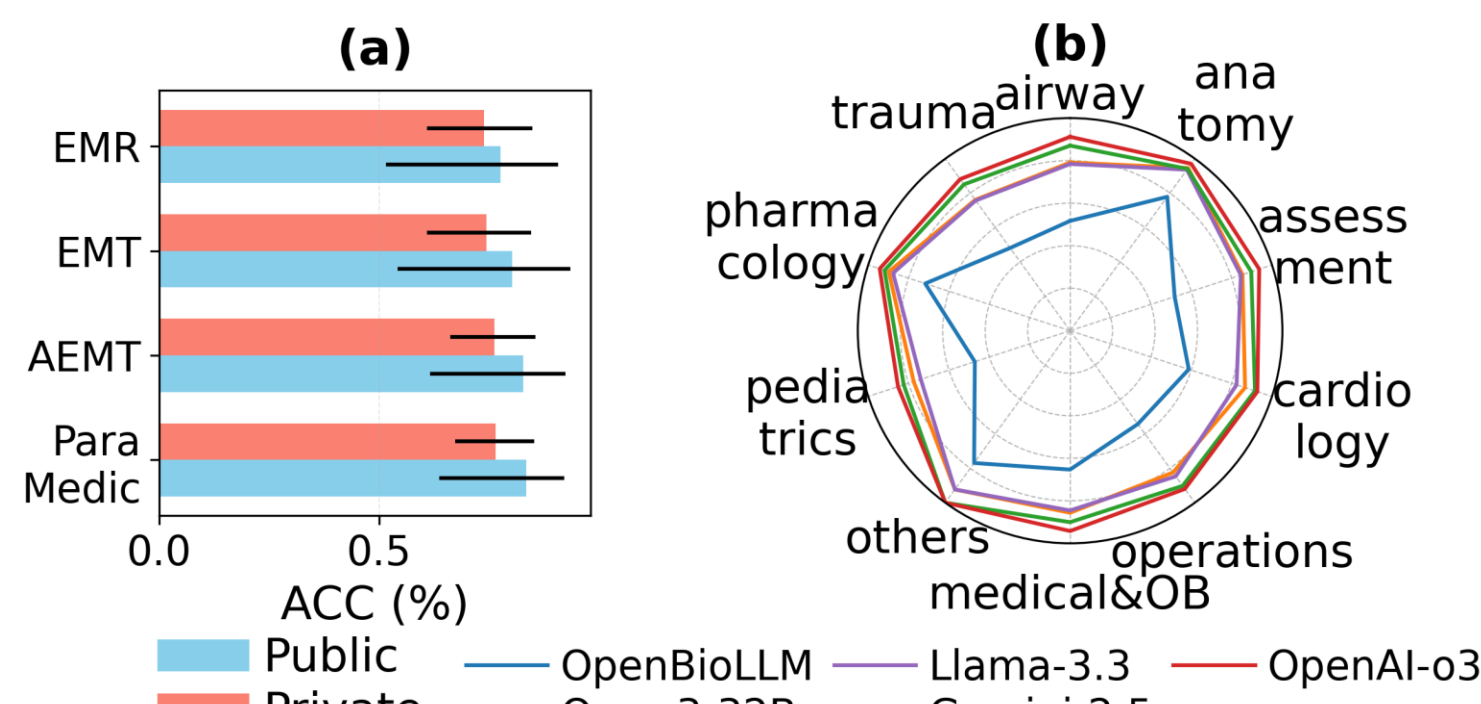
❖ ExpertRAG

$$\hat{A}_i = f^{\text{RAG}}(q_i, O_i, \mathcal{R}(q_i, \hat{s}_i), \hat{l}_i, \hat{s}_i).$$

- **Global (baseline):** Retrieve top M from KB and N from PR
- **Filter then Retrieve:** Filter documents in KB & PR using subject area \hat{s}_i , then retrieve top-M/N documents
- **Retrieve then Filter:** First retrieve 10×M/N candidates from KB/PR, then filter by \hat{s}_i to keep top-M/N docs

Results

❖ Where do SOTA LLMs shine or stumble across subject areas and certification?



❖ What is the Expertise Classification Performance?

Split	Model	Method	Subject Area miF	maF	Certification miF	maF
Public	Filter	LoRA	80.72	71.92	65.87	63.45
	Qwen3-4B	0-shot	55.43	51.61	45.77	30.49
	Qwen3-4B	4-shot	56.33	54.42	45.46	30.28
	Qwen3-4B	CoT	59.72	55.66	47.80	35.77
Private	Filter	LoRA	79.06	70.48	65.54	63.50
	Qwen3-4B	0-shot	42.93	31.73	44.12	25.01
	Qwen3-4B	4-shot	45.76	34.08	44.04	29.41
	Qwen3-4B	CoT	46.22	35.49	47.70	31.92

❖ How much does explicit expertise injected by Expert-CoT and ExpertRAG lift baseline accuracy?

Model	Prompt	Public Acc	Public F1	Private Acc	Private F1
No-RAG Baselines					
Qwen3-4B	0-shot	70.99	71.01	69.88	69.95
Qwen3-4B	CoT	72.35	73.09	70.58	72.02
RAG Baselines + CoT					
MedRAG	RAG on Med	74.31	74.41	71.12	73.33
i-MedRAG	Iterative RAG	77.96	78.00	74.02	76.35
Self-BioRAG	SelfRAG on Bio	55.71	58.84	45.72	49.67
Qwen3-4B	KB	76.49	76.07	75.02	76.53
Qwen3-4B	PR	73.02	73.96	70.54	72.38
Qwen3-4B	Global	78.12	79.17	75.46	76.87
RAG Baselines + Expert-CoT					
Qwen3-4B	KB	78.02	79.04	76.01	76.25
Qwen3-4B	PR	73.82	73.82	71.53	72.96
Qwen3-4B	Global	79.59	79.61	76.75	77.35
ExpertRAG-GT + CoT					
Qwen3-32B	0-shot	83.55	83.55	85.11	85.89
Qwen3-32B	4-shot	84.41	84.41	85.48	86.13
Qwen3-32B	32-shot	81.13	81.13	82.22	83.41
Qwen3-32B	64-shot	82.48	82.48	86.22	87.26
Qwen3-32B	CoT	84.96	84.97	88.78	90.13
Qwen3-32B	Expert-CoT	85.70	85.71	89.73	90.98
Qwen3-32B	Filter	85.57	85.60	89.50	91.20
ExpertRAG-GT + Expert-CoT					
OpenAI-o3	0-shot	92.39	92.39	-	-
Gemini-2.5	0-shot	89.36	89.36	-	-

❖ Can expertise-aware LLMs pass the NREMT standardized tests at different certification levels?

Model	Description	EMR				EMT				AEMT				Paramedic			
		Pass	Score	Acc	T	Pass	Score	Acc	T	Pass	Score	Acc	T	Pass	Score	Acc	T
Qwen3-4B	0-shot	X	809	64.18	33	X	940	74.07	33	X	940	71.64	33	X	940	71.72	35
	Expert-CoT	X	940	72.42	59	X	940	76.73	73	✓	1179	80.41	76	X	940	76.03	87
ExpertRAG-4B	FTR+Expert-CoT	✓	1218	84.21	66	✓	940	78.76	61	✓	940	77.31	69	✓	940	79.67	74
	RTF+Expert-CoT	✓	940	76.47	59	✓	1185	81.30	93	✓	1190	83.53	67	✓	940	77.61	86
Qwen3-32B	0-shot	✓	1207	82.65	22	✓	1140	81.63	23	✓	1280	85.92	26	✓	1163	81.58	29
	Expert-CoT	✓	1261	86.27	50	✓	1255	86.96	52	✓	1310	89.11	61	✓	1292	89.01	57
ExpertRAG-32B	FTR+Expert-CoT	✓	1350	92.22	75	✓	1292	89.01	76	✓	1215	84.60	86	✓	1228	83.93	125
	RTF+Expert-CoT	✓	1350	92.22	75	✓	1328	92.32	82	✓	1356	92.31	82	✓	1276	88.04	99

Acknowledgements

This work was supported by the award 70NANB21H029 from the U.S. Department of Commerce, National Institute of Standards and Technology (NIST), and a research grant from the Commonwealth Cyber Initiative (CCI).

