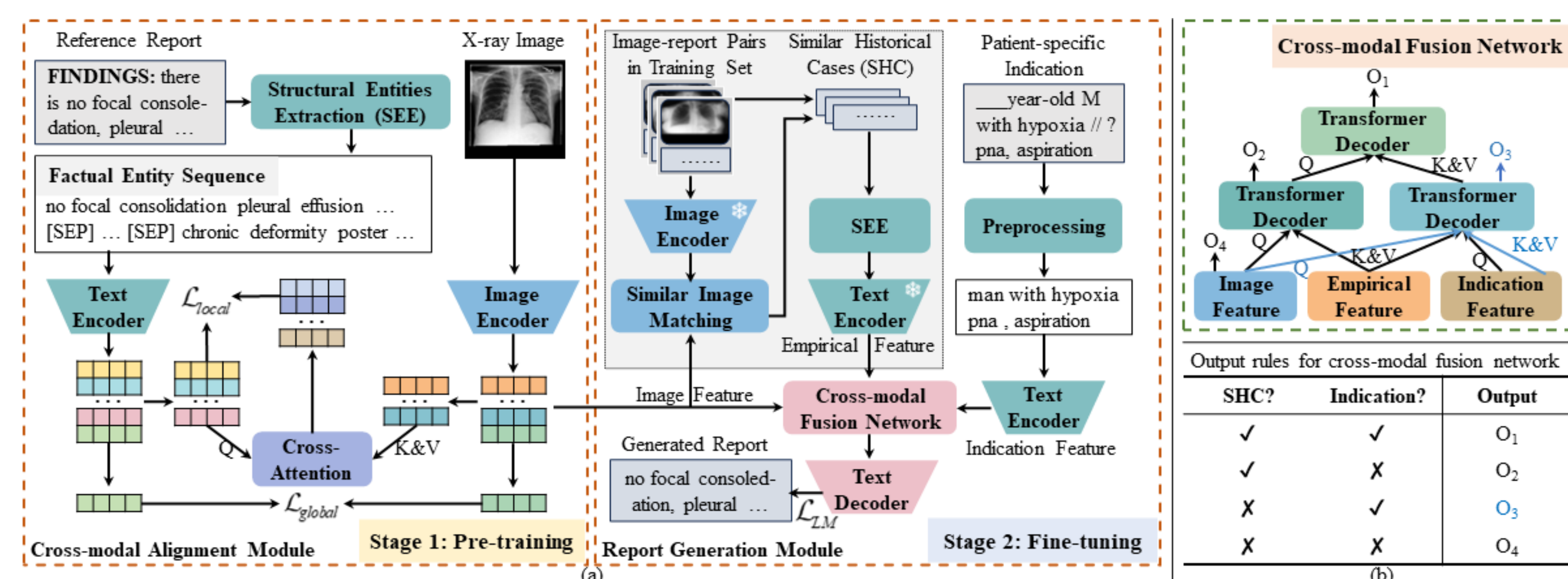


Kang Liu, Zhuoqi Ma\*, Xiaolu Kang, Zhushi Zhong, Zhicheng Jiao, Grayson Baird, Harrison Bai, Qiguang Miao\*

## Motivation

- ❑ **Background:** A radiology report comprises **presentation-style vocabulary**, which ensures clarity and organization, and **factual vocabulary**, which provides accurate and objective descriptions based on observable findings.
- ❑ To truly meet clinical needs, report generation processes should **incorporate patient-specific indications**, such as previous treatment history or responses to specific diagnostic requirements, which cannot be derived exclusively from medical images.
- ❑ Existing methods face challenges in effectively focusing on the cross-modal alignment between medical images and reports. This is attributed to assigning **equal weights to presentation-style elements** (e.g., sentence structure and grammar) **and factual vocabulary** (e.g., findings) in reports. Unfortunately, this limitation impacts their clinical efficacy.

## Method



We propose a **two-stage** method called **Structural Entities extraction and patient indications Incorporation (SEI)** for **generating chest X-ray report**.

- ✓ **Goal:** Given a chest X-ray, our model retrieves similar historical cases in a gradient-free manner and generates a draft report based on these cases for radiologists.
- ✓ **Stage 1: Pre-training for visual representation**
  - **Structural Entities Extraction (SEE):** We remove presentation-style vocabulary from reports using RadGraph outputs and process the noise in these outputs to form factual entity sequences, which are short sentences consisting exclusively of factual vocabulary.
  - **Cross-modal Alignment:** We align medical images with **factual entity sequence** through instance-level and token-level cross-modal semantic correspondences. This approach reduces the noise in the cross-modal alignment process, facilitating gradient-free retrieval of similar historical cases from the training set.
- ✓ **Stage 2: Fine-tuning for chest X-ray report generation**
  - **Patient-specific Indications:** This field is a string **representing the patient's examination purpose or symptoms**, which may occasionally be **absent**.
  - **Cross-modal Fusion Network:** This network effectively utilizes available indications and similar historical cases through output rules, even when some samples these elements. This enables the text decoder to attend to discriminative features of X-ray images, assimilate historical diagnostic information from similar cases, and understand the examination intention of patients.

## Experiments

- **Comparison of our SEI with SOTA approaches on the MIMIC-CXR dataset**

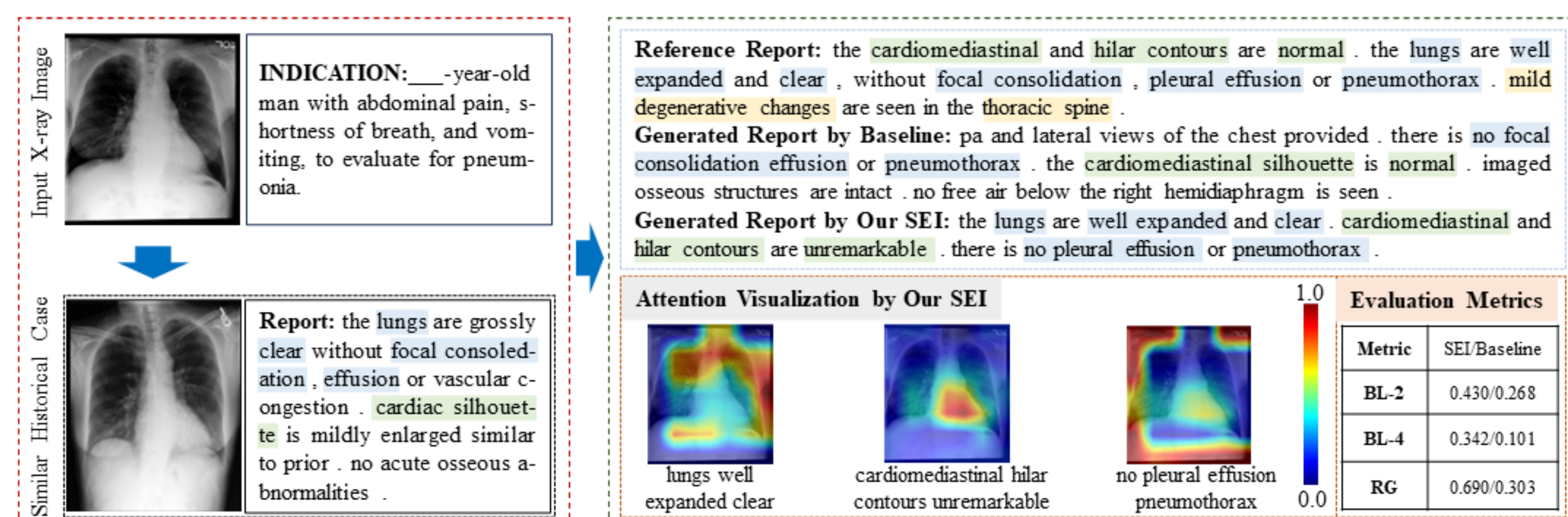
SEI- $n$  indicates our SEI generates reports by referencing  $n$  similar historical cases

Method	$M_{gt}$	NLG $\uparrow$				CE $\uparrow$		
		BL-2	BL-4	MTR	R_L	RG	CX5	CX14
R2Gen [5] (EMNLP'20)	100 $\uparrow$ <i>Cpl.</i>	0.218 0.209	0.103 0.097	0.137 0.135	0.264 0.266	0.207 0.211	0.340 0.339	0.340 0.338
R2GenCMN [4] (ACL'21)	100 $\uparrow$ <i>Cpl.</i>	0.218 0.198	0.106 0.090	0.142 0.133	0.278 0.268	0.220 0.223	0.461 0.464	0.278 0.393
GSKET [29] (MedIA'22)	80 $\uparrow$	0.228	0.115	-	0.284	-	-	0.371
CGPT2 [18] (ARTMED'23)	60 $\uparrow$ <i>Cpl.</i>	0.248 0.204	0.127 0.102	0.155 0.138	0.286 0.277	0.223 0.237	0.463 0.483	0.391 0.434
M2KT [28] (MedIA'23)	80 $\uparrow$ <i>Cpl.</i>	0.237 0.204	0.111 0.085	0.137 0.133	0.274 0.244	0.204 0.210	0.477 0.483	0.352 0.413
DCL [16] (CVPR'23)	90 $\uparrow$	-	0.109	0.150	0.284	-	-	0.373
RGRG [20] (CVPR'23)	<i>Cpl.</i> $\uparrow$	<b>0.249</b>	0.126	<b>0.168</b>	0.264	-	<b>0.547</b>	0.447
SEI-0 (ours)	60	<b>0.268</b>	0.146	0.164	0.300	<b>0.239</b>	0.505	0.437
	80	<b>0.250</b>	0.135	0.158	<b>0.300</b>	<b>0.250</b>	0.531	0.452
	90	<b>0.244</b>	0.131	0.156	0.299	<b>0.252</b>	0.536	0.455
	100	<b>0.240</b>	0.129	0.154	0.298	<b>0.252</b>	0.539	0.457
	<i>Cpl.</i>	0.231	0.123	0.150	<b>0.297</b>	<b>0.252</b>	0.541	0.457
SEI-1 (ours)	60	<b>0.268</b>	<b>0.148</b>	<b>0.167</b>	<b>0.301</b>	<b>0.236</b>	<b>0.509</b>	<b>0.445</b>
	80	<b>0.257</b>	0.140	<b>0.162</b>	<b>0.300</b>	0.247	<b>0.535</b>	<b>0.457</b>
	90	<b>0.251</b>	<b>0.137</b>	<b>0.160</b>	<b>0.300</b>	0.248	<b>0.539</b>	<b>0.459</b>
	100	<b>0.247</b>	<b>0.135</b>	<b>0.158</b>	<b>0.299</b>	0.249	<b>0.542</b>	<b>0.460</b>
	<i>Cpl.</i>	<b>0.238</b>	<b>0.128</b>	0.154	0.296	0.249	<b>0.545</b>	<b>0.460</b>

- **Ablation study on MIMIC-CXR**

Settings	Model	NLG $\uparrow$				CE $\uparrow$		
		BL-2	BL-4	MTR	R_L	RG	CX5	CX14
(a)	Base (R2Gen [5])	0.209	0.097	0.135	0.266	0.211	0.339	0.338
(b)	(a)+cross-modal module	0.206	0.098	0.138	0.277	0.234	0.513	0.431
(c)	SEI-1 w/o indications	0.228	0.109	0.148	0.279	0.241	0.542	<b>0.474</b>
(d)	SEI-1 w/o SHC (SEI-0)	0.231	0.123	0.150	<b>0.297</b>	<b>0.252</b>	0.541	0.457
(e)	<b>SEI-1</b>	<b>0.238</b>	<b>0.128</b>	<b>0.154</b>	0.296	0.249	<b>0.545</b>	0.460

- **Qualitative analysis on MIMIC-CXR**



- **Conclusion:**
  - ✓ Using **factual entity sequences for alignment** proves to be an effective strategy.
  - ✓ **Similar historical cases** provide valuable empirical insights for report generation.
  - ✓ Integrating **patient-specific indications** into the report generation process significantly enhances its performance.

Welcome to our Github homepage for source codes  
<https://github.com/mk-runner/SEI>

Email: kangliu422@gmail.com  
WeChat: kangliu422

