

# Using Decoder-Based Distillation for Enhancing Multilingual Clinical Case Report Summarization

Nicolay Rusnachenko<sup>1</sup> Xiaoxiao Liu<sup>1</sup> Jian Chang<sup>1</sup> Jian Jun Zhang<sup>1</sup>

<sup>1</sup> Centre for Applied Creative Technologies (CFACT+), Faculty of Media and Communications, Bournemouth, United Kingdom

## Key Takeaway

- Distillation framework with *role-based dialogue modeling* notation for **Student Models** with **Clinical Key Info** derived from reports via **Teacher Model**
  - Exploit of **System**, **User**, and **Assistant** roles which are commonly supported by instruction-tuned models
- We experiment with technique adaptation in clinical case report summarization task for **Qwen-2.5** models family
  - extracting clinical key information from *teacher model* (72B params) and using this information in tuning of small-scaled student model (0.5B params) **results in 2.4%-4% on MultiClinSum<sup>small</sup>** while at evaluation stage.

YouTube WATCH ONLINE



## Related Works

- Xiaoxiao Liu and et. al.  
Enhancing medical dialogue summarization: A medextract distillation framework.  
pages 6466–6473, 2024.
- M Rodríguez-Ortega and et. al.  
Overview of multiclinsum task at bioasq 2025: evaluation of clinical case summarization strategies for multiple languages: data, evaluation, resources and results.

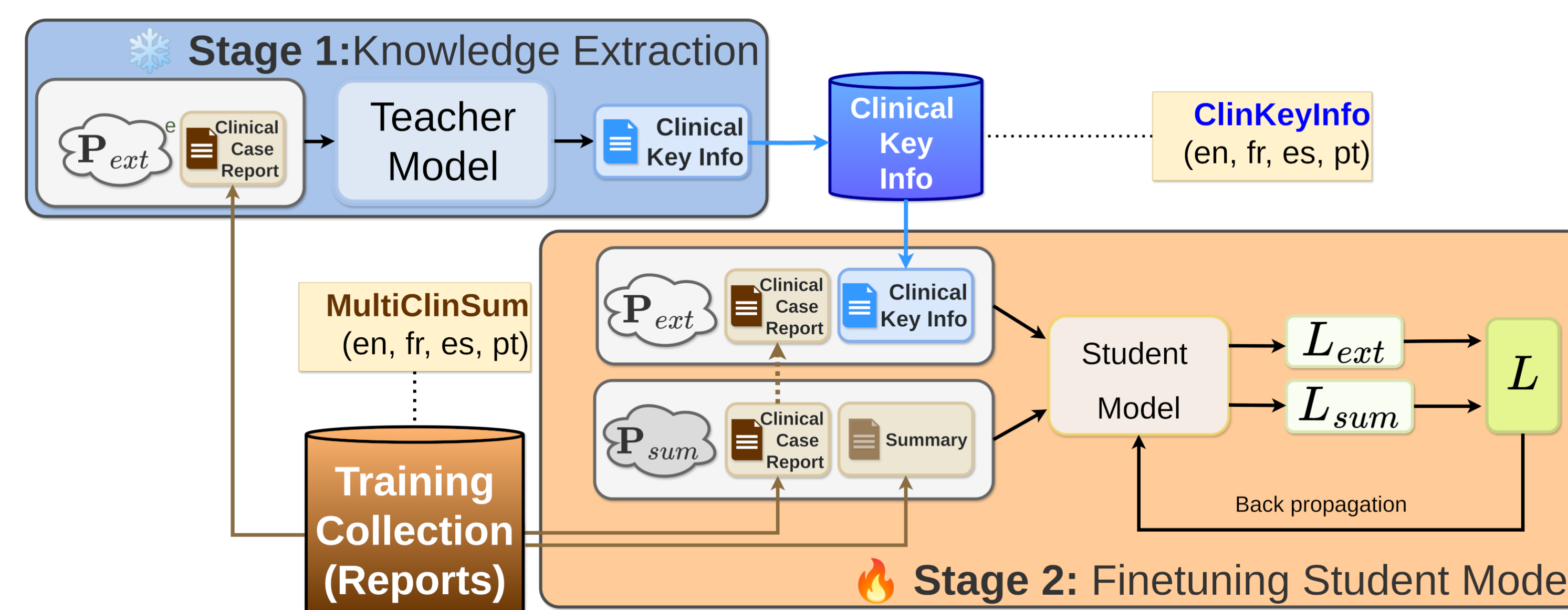
## Implementation



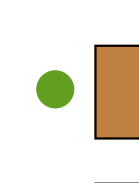
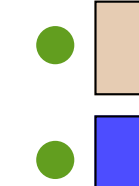




nicolay-r/distil-tuning-llm




## Two-stage distillation for role-based input systems




### Definitions:

-  – Clinical Case Report from Reports Collection
-  – Clinical Case Report Summary (from Reports Collection)
-  – Clinical Key information from Clinical Key Information collection
  - extracted from  on Stage 1
- $P_{ext}$  – ext. prompt: «Extract the key information from clinical text: »
- $P_{sum}$  – summarization prompt: «Summarize clinical text: »
- (S, U, A) – role-based dialogue modelling notation.

### Stage 1

Given: [ $P_{ext}$ , Dataset of ]




For each :




- Infer result from Teacher Model:
- (S:  $P_{ext}$ , U: , A:  $\emptyset$ )

Result: dataset of 

### Stage 2

Given:

- $P_{ext}$  – extraction prompt;
- $P_{sum}$  – summarization prompt;
- Dataset of (, , )

For each (, , ):

- Extraction Supervision ( $L_{ext}$ ):

$$(S: P_{ext}, U: \text{Clinical Case Report}, A: \text{Clinical Key Info})$$

$$L_{ext} = \sum_{t=i_{start}}^{T_e} l(y_e, \hat{y}_e, t, x_e)$$

- Summarization Supervision ( $L_{sum}$ ):

$$(S: P_{sum}, U: \text{Clinical Case Report}, A: \text{Clinical Case Report Summary})$$

$$L_{sum} = \sum_{t=i_{start}}^{T_s} l(y_s, \hat{y}_s, t, x_s)$$



- Combined loss ( $L$ ) with  $\gamma = 0.8$ :

$$L = \gamma L_{sum} + (1 - \gamma) L_{ext}$$


Result: Fine-tuned Student Model

### Dataset I

- MultiClinSum<sup>small</sup>** – 592 reports of each language from original MultiClinSum:
- ClinKeyInfo<sup>small</sup>** – augm. of reports of MultiClinSum<sup>small</sup> with  obtained from Stage 1




Lng	#	# Chars in 			# Chars in 		
		Avg	Min	Max	Avg	Min	Max
EN	592	3785.4	719	34071	725.3	90	3883
ES	592	4056.1	825	17602	792.6	125	4161
FR	592	4783.2	827	37138	832.1	121	4542
PT	592	4096.0	793	37351	809.5	116	28227


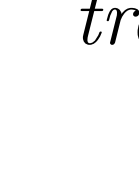



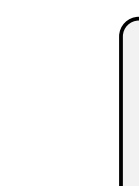
### Dataset II

Lng	#	# Chars in 		
		Avg	Min	Max
EN	592	2971.4	1088	6597
ES	592	2929.5	392	6040
FR	592	2873.3	879	5472
PT	592	2871.9	911	7961

- We use Qwen-2.5-72B-instruct.

## Experimental Setup

Our setup based on content from **MultiClinSum<sup>small</sup>** (, ) and **ClinKeyInfo<sup>small</sup>** () for subsets:

Subset	#	# Chars Avg Range			
train		1892	2435.6	719-2560	
		1892	490.6	90-512	
		1892	511.9	392-512	
valid		20	2560.0	2560-2560	
		20	510.2	486-512	
		20	512.2	512-512	

We use Qwen2.5-0.5B as a student model for Stage 2 to prepare:



## Non-official Results

### Models

- baseline Qwen2.5-0.5B
- † – finetuned Qwen2.5-0.5B<sub>standard</sub>
- ‡ – distil-tuned Qwen2.5-0.5B<sub>distil</sub>

### Test Set

The 456 of original reports were used (19% of MultiClinSum<sup>small</sup>)

L	BERTScore			ROUGE		
	P	R	F1	R-1	R-2	R-L
E •	78.59	<b>82.91</b>	80.62	32.49	11.88	20.97
N †	80.94	82.13	81.47	37.49	15.26	25.48
	‡	<b>81.80</b>	81.67	<b>81.69</b>	<b>38.30</b>	<b>15.57</b>
E •	80.26	<b>84.64</b>	82.35	33.66	13.33	20.69
S †	84.07	83.62	<b>83.80</b>	<b>40.50</b>	<b>17.14</b>	<b>26.74</b>
	‡	<b>84.10</b>	83.48	83.76	40.26	16.72
F •	81.16	<b>84.36</b>	82.69	34.45	13.80	20.23
R †	84.05	83.80	<b>83.88</b>	<b>39.67</b>	<b>17.00</b>	<b>24.94</b>
	‡	<b>84.34</b>	83.10	83.68	38.95	16.25
P •	80.65	<b>83.53</b>	82.02	30.81	11.44	19.53
T †	83.21	83.29	83.19	37.51	<b>15.02</b>	<b>24.30</b>
	‡	<b>83.42</b>	83.10	<b>83.22</b>	14.92	24.19

## Official Results

Model: Qwen2.5-0.5B<sub>distil</sub>

Lng	BERTScore			ROUGE		
	P	R	F1	P	R	F1
EN	85.54	85.70	85.59	27.53	27.53	25.87
ES	72.42	73.47	72.88	26.06	29.03	25.87
FR	72.48	73.96	73.15	24.15	28.90	24.66
PT	72.39	73.20	72.73	24.95	27.05	24.40