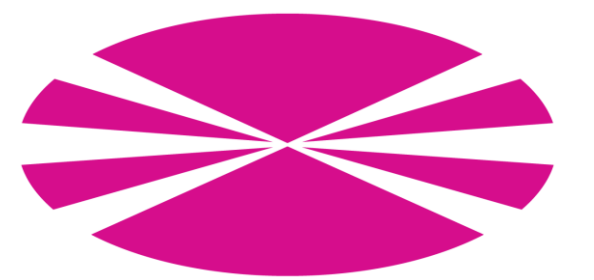


# LyS@SemEval-2025 Task 8: Zero-Shot Code Generation for Tabular QA



Adrian Gude, Roi Santos-Ríos, Francisco Prado-Valiño, Ana Ezquerro AND Jesús Vilares

Centro de Investigación en TIC (CITIC), Universidade da Coruña, Spain



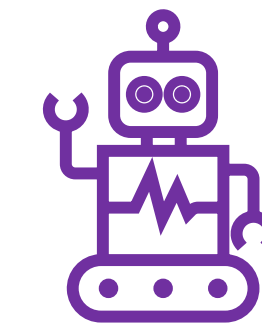
## TABULAR QA AS CODE GENERATION

Leverage powerful LLMs to *generate* executable code.

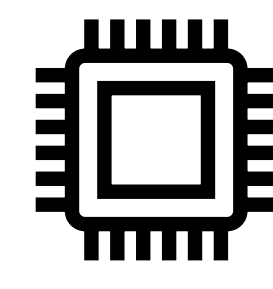
- Column selection to identify relevant variables and data types (Herzig et al., 2020).
- Code fixer module: capture and integrate runtime errors in the prompt.

**Result:** Guided prompt-based code generation (no training/finetuning).

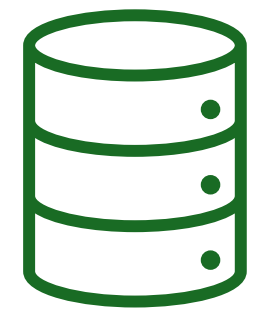
Key components of the system



LLM



Compiler



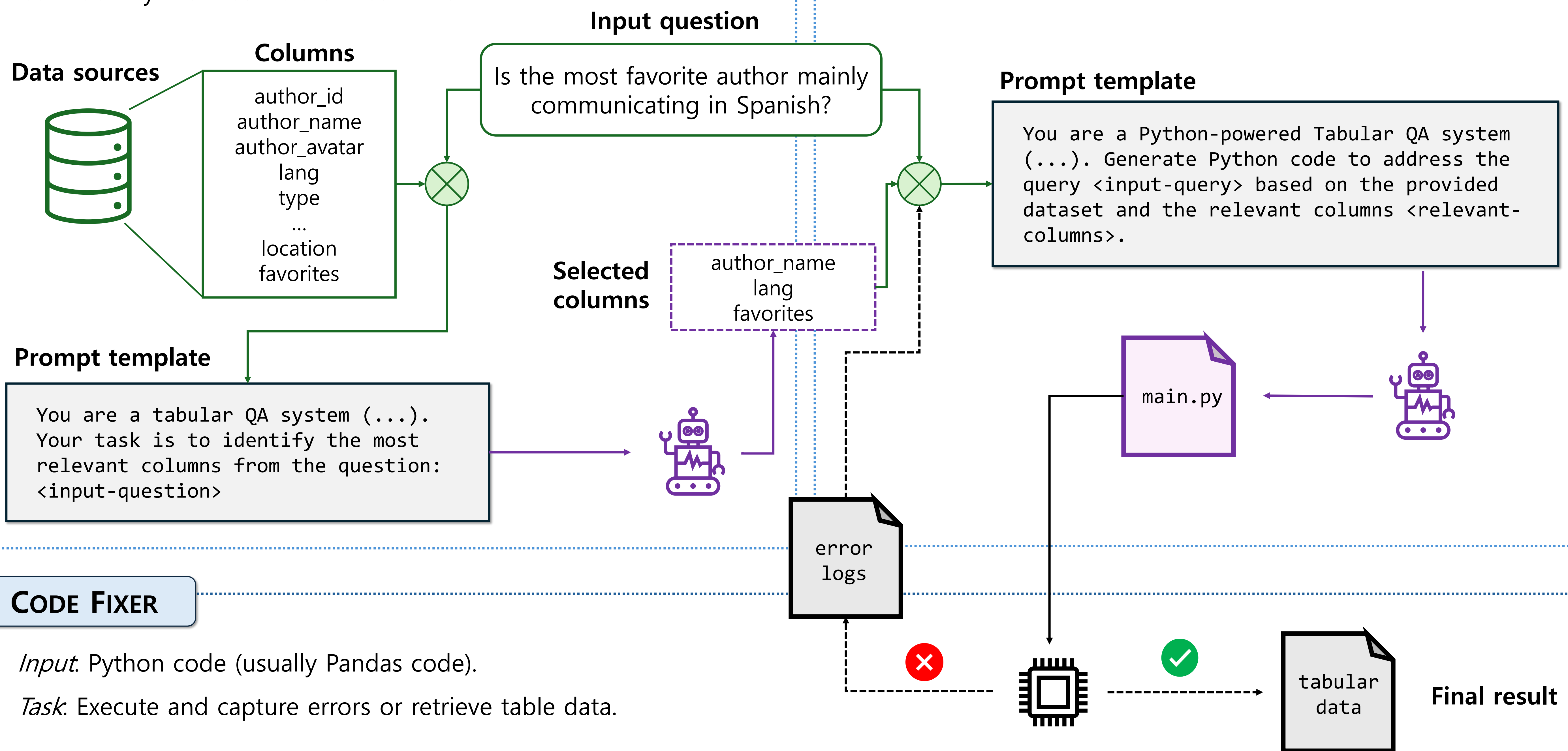
Data source

## COLUMN SELECTOR

- Input:* question (natural language) + list of columns.
- Task:* Identify the most relevant columns.

## ANSWER GENERATOR

- Input:* question (natural language) + selected columns.
- Task:* Generate Python code.



## CODE FIXER

- Input:* Python code (usually Pandas code).
- Task:* Execute and capture errors or retrieve table data.

## EXPERIMENTS AND RESULTS

- Better performance than the baseline  $\beta$  (S1: 27.00; S2: 26.00).
- One of the best systems in the development phase (top 5).
- Ablation study with different LLMs.
  - Qwen-2.5-Coder (7B, 32B), Codestral (22B), Mistral (7B).
- Ablation study with different configurations:
  - Enhanced Column Selection (ECS)
  - Answer Generator (AG).
  - Code Fixer (CF).
- Final phase:**
  - General category: 32nd (S1) and 31st (S2) out of 49 participants.
  - Open models category: 23rd (S1) and 21st (S2) out of 35 participants.

Ablation study with Qwen-2.5-Coder<sup>32B</sup> on the dev set

		boolean	category	number	list[category]	list[number]	$\mu$
S1	AG	81.25	78.12	75.00	65.62	70.31	74.06
	AG+CS	82.81	78.12	78.12	68.75	79.69	77.50
	AG+ECS+CF	<b>89.06</b>	<b>85.94</b>	<b>85.94</b>	<b>78.12</b>	<b>85.94</b>	<b>85.00</b>
S2	AG	84.37	<b>89.06</b>	85.94	75.00	75.00	81.87
	AG+CS	84.37	<b>89.06</b>	<b>90.62</b>	73.44	<b>79.69</b>	83.44
	AG+ECS+CF	<b>89.06</b>	<b>89.06</b>	<b>90.62</b>	<b>76.56</b>	78.12	<b>84.69</b>

Ablation study of different LLMs on the dev set (only Answer Generator)

		boolean	category	number	list[category]	list[number]	$\mu$	$\beta$
S1	Qwen-2.5-Coder <sup>7B</sup>	67.19	68.75	75.00	3.12	3.12	43.44	27.00
	Mistral <sup>7B</sup>	51.56	59.37	73.44	35.94	34.37	50.94	
	Codestral <sup>22B</sup>	73.44	<b>82.81</b>	48.44	48.44	48.44	67.19	
	Qwen-2.5-Coder <sup>32B</sup>	<b>81.25</b>	78.12	75.00	<b>65.62</b>	<b>70.31</b>	<b>74.06</b>	
S2	Qwen-2.5-Coder <sup>7B</sup>	81.25	84.37	85.93	6.25	1.56	51.87	26.00
	Mistral <sup>7B</sup>	46.87	56.25	65.62	32.81	25.00	45.31	
	Codestral <sup>22B</sup>	71.87	<b>89.06</b>	84.37	53.12	60.94	71.87	
	Qwen-2.5-Coder <sup>32B</sup>	<b>84.37</b>	<b>89.06</b>	<b>85.94</b>	<b>75.00</b>	<b>75.00</b>	<b>81.87</b>	

Average performance ( $\mu$ ) and baseline ( $\beta$ ) included. Best performance in bold.

