

# A Benchmark Dataset for Identifying Machine-Generated Scientific Papers in the LLM Era.



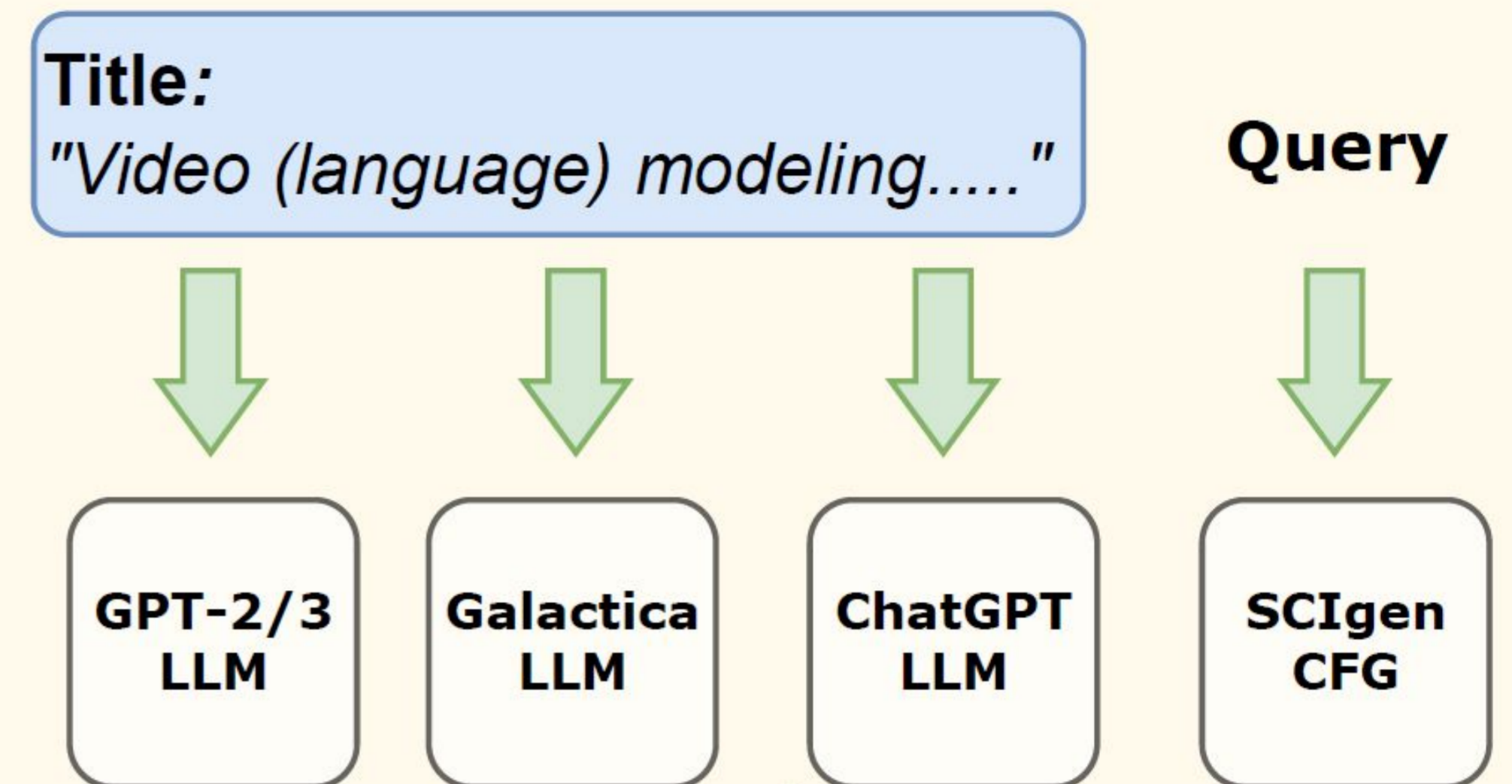
Edoardo Mosca, Mohamed Hesham I. Abdalla, Paolo Basso,  
Margherita Musumeci, Georg Groh

## Contributions

- **Benchmark dataset** comprising **real (human-written)** and **fake (machine-generated) scientific documents**. Each contains an **abstract**, **introduction**, and **conclusion**.
- Evaluation of **four different classifiers** to **determine the authorship** (real or fake) of the documents.
- Analysis and expla of **classifiers' generalization abilities** by evaluating their performance on both **in-domain** and **out-of-domain settings**.

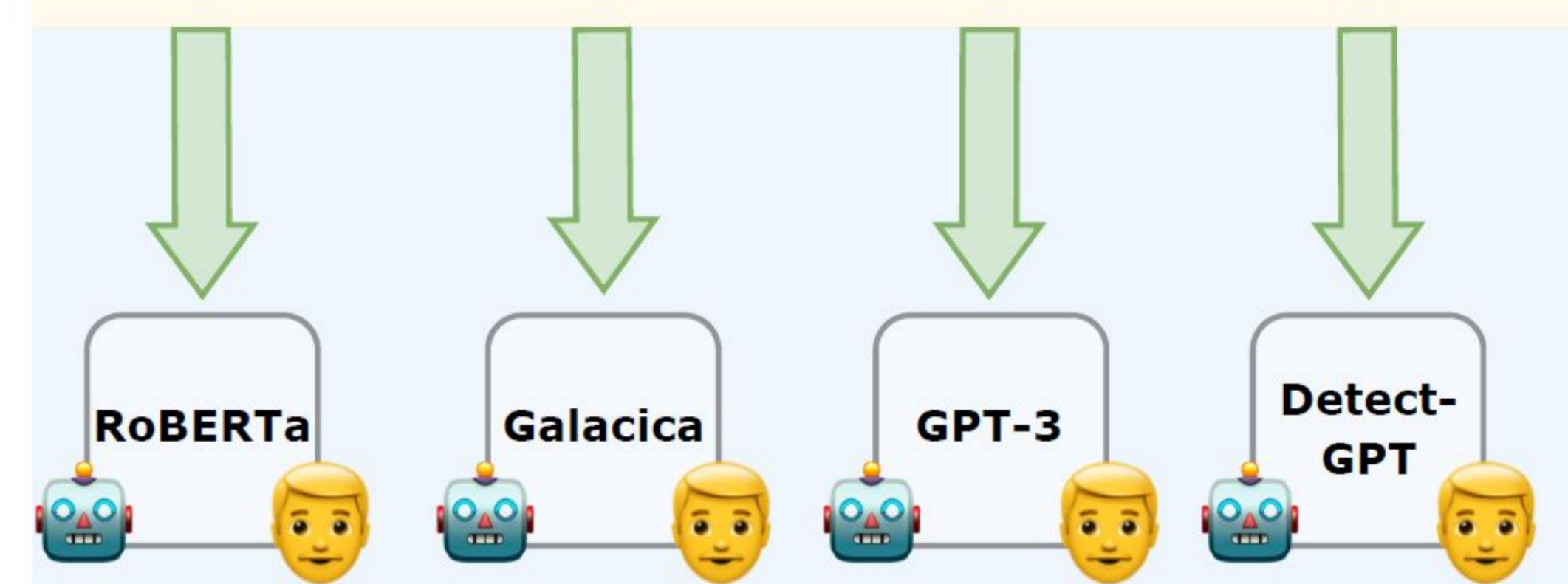
Can we distinguish facts from fiction?

Benchmark Generation

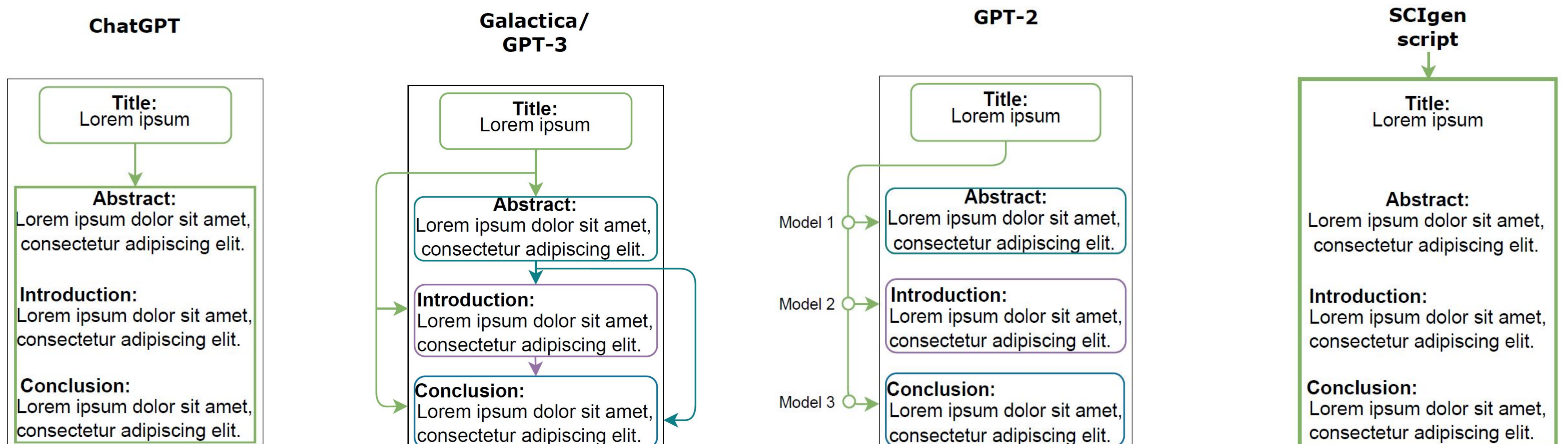


**Abstract:** "Advances in video modeling..  
**Introduction:** "Video data is a growing..  
**Conclusion:** "In our work, we tested the.."

Detection



## Generation Sequence



## Overview of the Dataset

Source	Quantity	Tokens
arXiv parsing 1 (real)	12k	13.40M
arXiv parsing 2 (real)	4k	3.20M
SCIgen (fake)	3k	1.80M
GPT-2 (fake)	3k	2.90M
Galactica (fake)	3k	2.00M
ChatGPT (fake)	3k	1.20M
GPT-3 (fake)	1k	0.50M
<b>Total real (extraction)</b>	<b>16k</b>	<b>16.60M</b>
<b>Total fake (generators)</b>	<b>13k</b>	<b>8.40M</b>
<b>Total</b>	<b>29k</b>	<b>25M</b>

## Experimental Results

Model	Train Dataset	TEST	OOD-GPT3	OOD-REAL	TECG
GPT-3 (our)	TRAIN-SUB	<b>99.96%</b>	25.9%	99.07%	<b>100%</b>
Galactica (our)	TRAIN	98.3%	24.6%	95.8%	83%
Galactica (our)	TRAIN+GPT3	98.5%	70%	92.1%	87.2%
Galactica (our)	TRAIN-CG	95%	11.1%	96.9%	42%
RoBERTa (our)	TRAIN	86%	23%	76%	<b>100%</b>
RoBERTa (our)	TRAIN+GPT3	68%	<b>100%</b>	36%	63%
RoBERTa (our)	TRAIN-CG	75%	32%	58%	88%
DetectGPT	-	61.5%	0%	<b>99.92%</b>	68.7%

■ Indicates out-of-domain experiments

## Takeaways

- Detection baselines can be really good, but sometimes **struggle with out-of-domain data**.
- No good **open-source detectors** available to test against.
- Future research should **include more paper sections**.
- **Human-Machine hybrids** are also a must for future research.



Data  
Models  
Code

