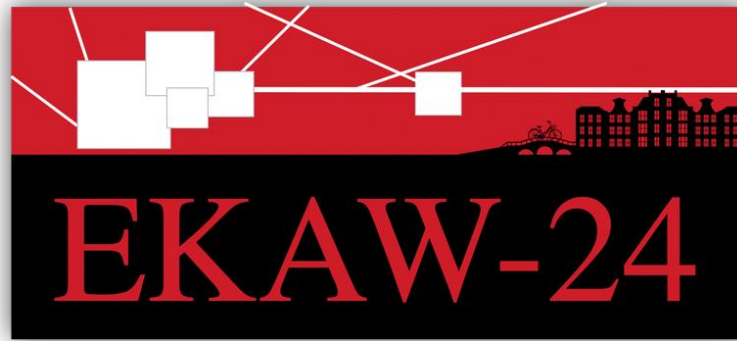


# Named Entity Recognition in Historical Italian: The Case of Giacomo Leopardi's Zibaldone

Cristian Santini, Laura Melosi, Emanuele Frontoni  
*University of Macerata*



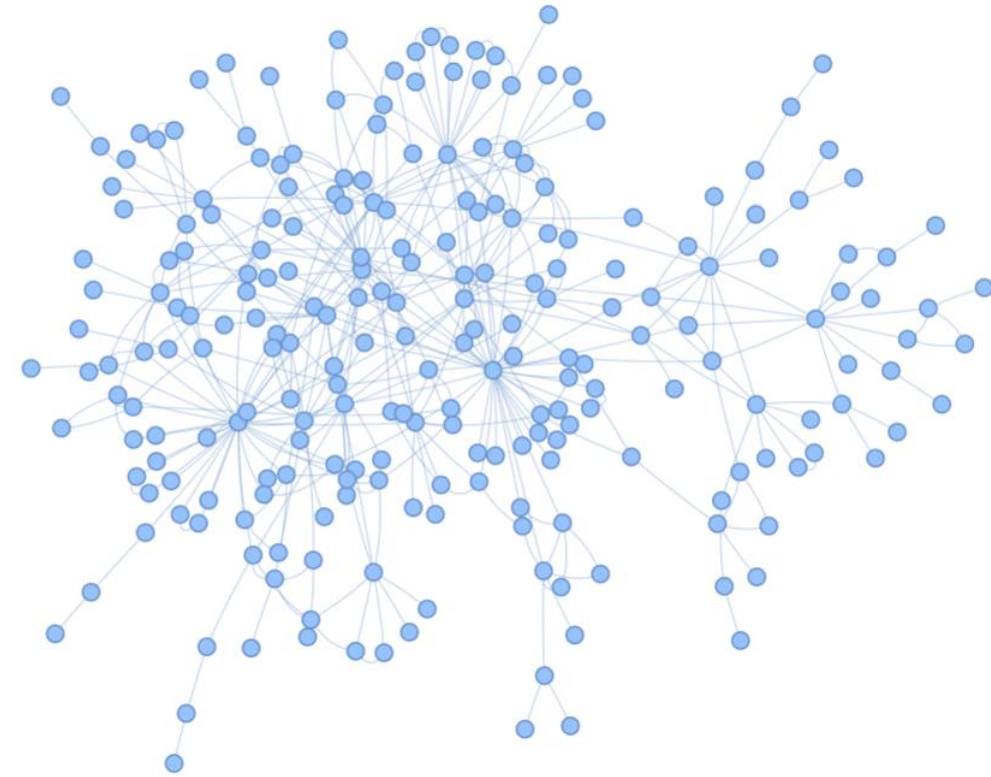
# Context

- In 2017, the idea to catalogue and digitize the autograph manuscripts of Giacomo Leopardi was launched
- Since 2018, the University of Macerata has joined the project, renamed 'Biblioteca Digitale Leopardiana' [1]
- The goal is to launch a digital platform aimed at collecting images and metadata of the documents, as well as at enhancing Leopardi's texts and studies



# Leopardi's Digital Heritage

- Currently, more than 15,000 *facsimilia* of Giacomo Leopardi are digitized and publicly available on the web
- Efforts have been done to provide semi-diplomatic editions of his works [2] [3]
- Named Entity Recognition (NER) techniques are therefore needed to model facts expressed in natural language into a tailored Knowledge Graph



[2] Gioele Marozzi (Ed.), *Giacomo Leopardi*. Cambridge Digital Library. <https://cudl.lib.cam.ac.uk/collections/leopardi/1>

[3] Silvia Stoyanova and Ben Johnston (Ed.), *Giacomo Leopardi's Zibaldone di pensieri: a digital research platform*. <https://digitalzibaldone.net/>

# Background: LLMs and Historical Texts

- Named Entity Recognition with ChatGPT on literary texts (ita,tr) [4] [5]
  - ✓ Applicable in absence of training data
  - ✗ No quantitative evaluation
- Benchmarking LLMs on historical multilingual press articles [6]
  - ✗ Vulnerable to noise and digitization errors
  - ✗ Challenges of mitigating hallucinations and controlling output
  - ✗ Worse performance than fine-tuned models

[4] S. Spina, Biscari Epistolography. From Archive to the Website., DigItalia 18 (2023) 245–259. <https://digitalia.cultura.gov.it/article/view/3010>.

[5] F. Aladağ, The Potential of GPT in Ottoman Studies: Computational Analysis of Evliya Celebi's Travelogue with NLP and Text Mining and Digital Edition with TEI, CULTURE 5 (2023).

[6] González-Gallardo, C.-E., Hanh, T. T. H., Hamdi, A., & Doucet, A. (2024). Leveraging Open Large Language Models for Historical Named Entity Recognition. The 28th International Conference on Theory and Practice of Digital Libraries. <https://univ-rochelle.hal.science/hal-04662000>

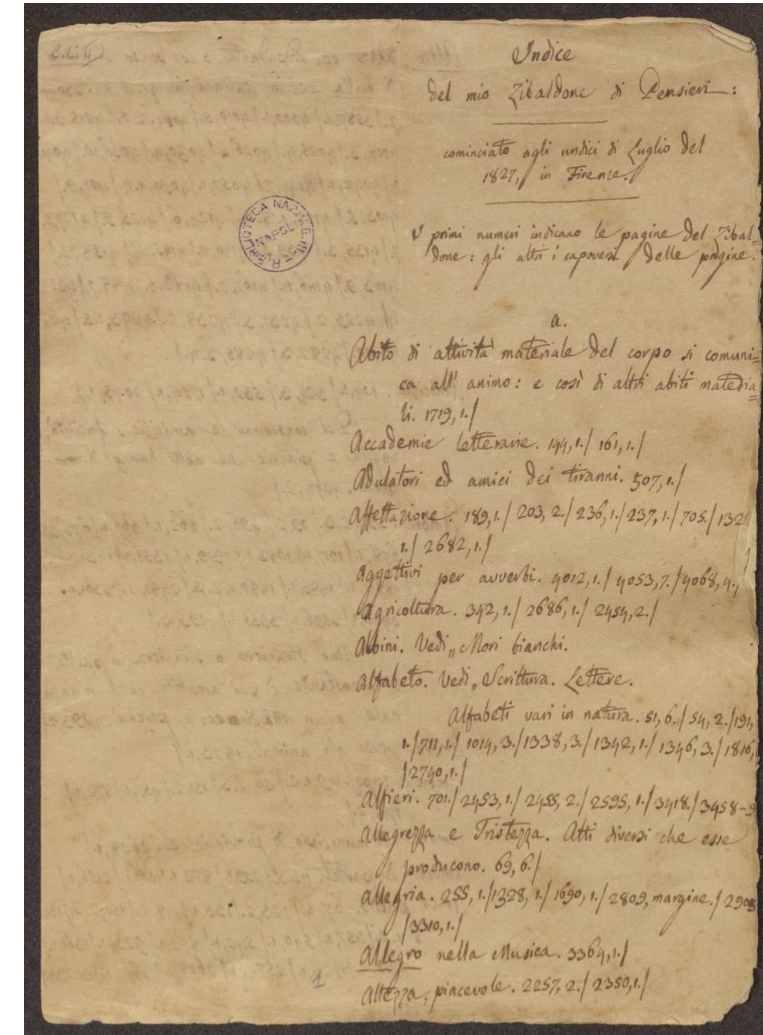
# Contribution

1. Propose a new dataset for Named Entity Recognition (NER) and Entity Linking (EL) extracted from Leopardi's *Zibaldone* (1817-1832)
2. Understand how efficient LLMs and fine-tuned models can be in recognizing references to named entities in Nineteenth century Italian humanistic texts
3. Give insights on the challenges of NER models and LLMs with Italian historical and humanistic texts



# DigitalZibaldone

- DigitalZibaldone [3] is a digital edition of Leopardi's diary made by scholars in Princeton
- More than 4,000 encoded notes containing observations on a wide-range of topics (linguistics, history, philosophy, philology, politics, etc.)
- Over 10,000 references to people, places and literary works linked to Wikidata and VIAF and encoded in HTML



# Markup of DigitalZibaldone

```
<p xmlns="http://www.w3.org/1999/xhtml" xmlns:tei="http://www.tei-c.org/ns/1.0">
<div class="node" id="p2721_1"><div class="nodemeta" id="p2721_1_meta"></div><b class="para_num">[2721,1]</b><u>NBSP</u> Anche il <a
class="person" href="https://digitalzibaldone.net/node/Q518160">Gelli</a> confessava (ap. <a class="person"
href="https://digitalzibaldone.net/node/Q3769747">Perticari</a>
<a href="/node/viaf34613848">Degli Scritt. del Trecento</a> l. 2. c. 13. p. 183.) che la lingua toscana
non era stata applicata alle scienze. (24. Maggio 1823.).</div></p>
```

# Data Sampling

- **Evaluation dataset:** 260 notes sampled from pages 2700-3000
  - All written in 1823 → most productive year for Leopardi
- **Training dataset:** 688 notes sampled in page ranges 1000-2000 and 3001 – 4000 containing at least one entity
- Filtering documents in the training set longer than 350 tokens → keep texts with similar length



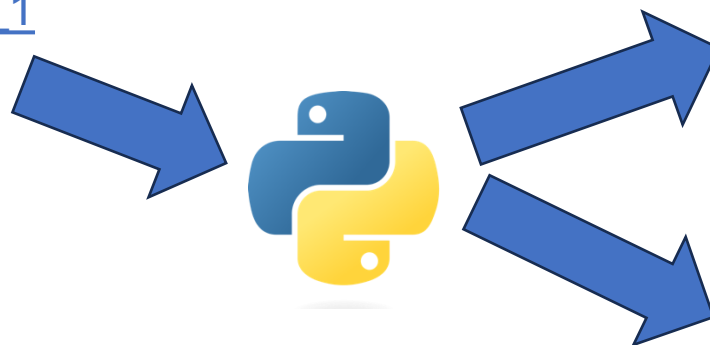
# Data Pre-processing

« 2718,1

2721,2»

**[2721,1]** Anche il Gelli confessava (ap. Perticari Degli Scritt. del Trecento I. 2. c. 13. p. 183.) che la lingua toscana non era stata applicata alle scienze. (24. Maggio 1823.).

[https://digitalzibaldone.net/node/p2721\\_1](https://digitalzibaldone.net/node/p2721_1)



doc_id	text
p2721_1	Anche il Gelli confessava (ap. Perticari [...])



doc_id	surface	start_pos	end_pos	type	identifier
p2721_1	Gelli	9	14	PER	Q518160
p2721_1	Perticari	31	40	PER	Q3769747
p2721_1	Degli Scritt. del Trecento	41	67	WORK	viaf34613848



# Dataset Statistics

## Training

Statistic	Value
Number of documents	688
Number of annotations	2,135
Person Annotations	1,093
Location Annotations	407
Work Annotations	635

## Evaluation

Statistic	Value
Number of documents	260
Number of annotations	764
Person Annotations	492
Location Annotations	61
Work Annotations	211

# Models Tested

- **Instruction-tuned LLaMa3.1 8B [7]:** enhanced in-context learning capabilities for zero-shot NER
- **GLiNER [8]:** BERT-based NER model trained on Italian texts to jointly learn entity and class representations
  - **Advantages:** can be evaluated both fine-tuned and in zero-shot setting

[7] Dubey, A., Jauhri, A., Pandey, A., Kadian, A., ... Zhao, Z. (2024). The Llama 3 Herd of Models. arXiv. (arXiv:2407.21783).

<https://doi.org/10.48550/arXiv.2407.21783>

[8] Zaratiana, U., Tomeh, N., Holat, P., & Charnois, T. (2023). GLiNER: Generalist Model for Named Entity Recognition using Bidirectional Transformer (arXiv:2311.08526). arXiv. <https://doi.org/10.48550/arXiv.2311.08526>

# LLaMa Prompts

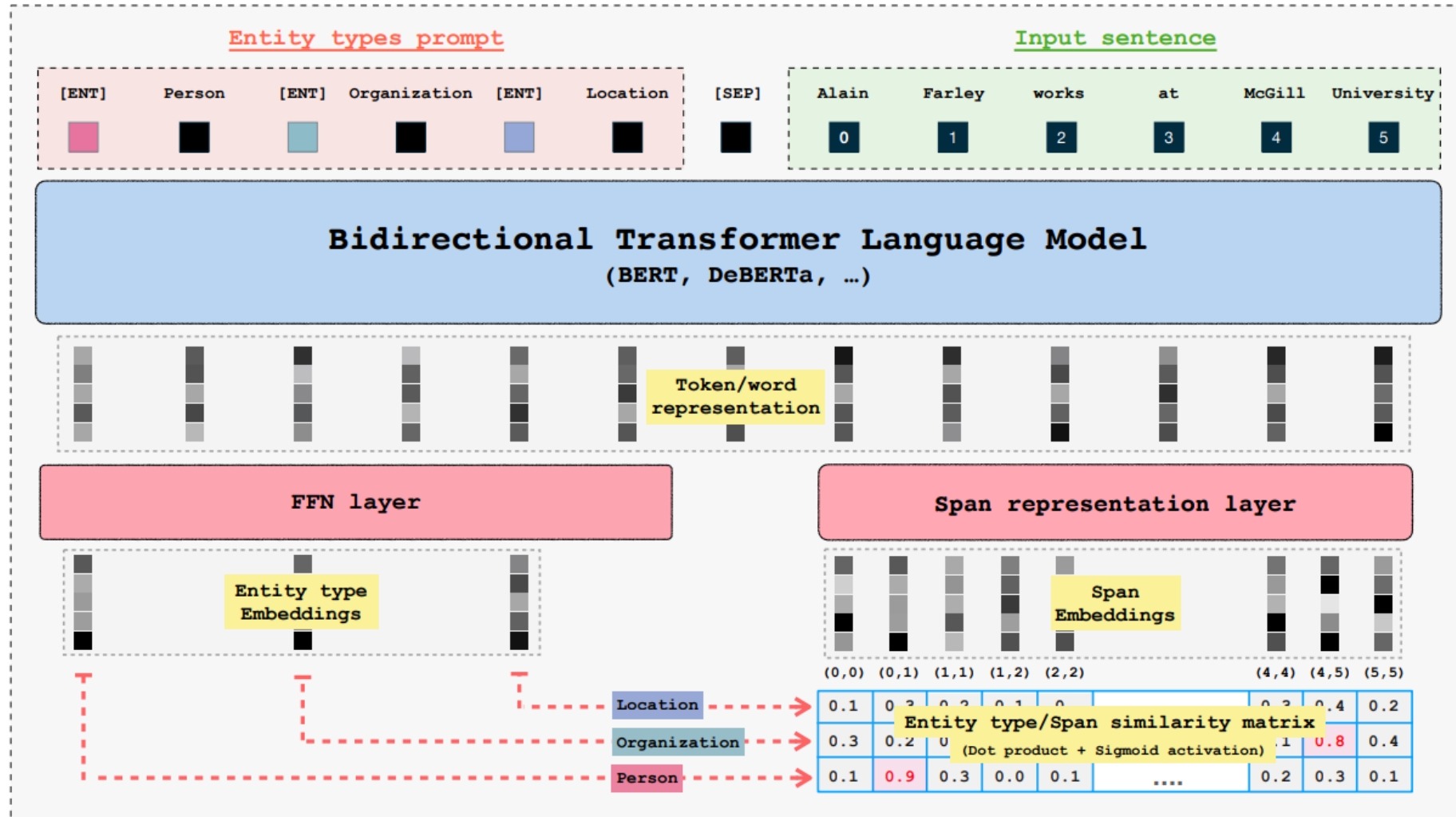
## GENERATIVE PROMPT

Here's the annotated text: «The first example of `<per>`Dante`</per>`'s `<work>`Divine Comedy`</work>` is owned by `<loc>`Biblioteca Passerini-Landi`</loc>` in `<loc>`Piacenza`</loc>`.»

## EXTRACTIVE PROMPT

Here is the list of entities:  
`<per>`Dante`</per>`,  
`<work>`Divine Comedy`</work>`,  
`<loc>`Biblioteca Passerini-Landi`</loc>`,  
`<loc>`Piacenza`</loc>`.

# GLiNER



# Evaluation Setting

Model	Setting 1	Setting 2
LLaMa3.1-instruct-8B	<b>Generative Prompt:</b> the LLM rewrites the annotated text based on a set of target classes	<b>Extractive Prompt:</b> the LLM extracts a list of entities from the text based on a set of target classes
GliNER_ita_base	<b>Zero-shot:</b> the model tries to identify unseen classes having been trained only on general texts	<b>Fine-tuned:</b> the model is pre-trained and fine-tuned on the training portion of the Zibaldone



# Evaluation Metrics

- NER results computed using both *exact* and *fuzzy* matching
  - Allows to assess which class is more affected by boundary detection errors
- Precision, Recall and F1 computed for each algorithm
  - Micro-averaged across all classes
  - Computed separately for each class and macro-averaged

# Micro-averaged Results

	Exact			Fuzzy		
	Precision	Recall	F1	Precision	Recall	F1
LLaMa3.1-8B (generative)	22,48	48,42	30,71	24,73	53,27	33,78
LLaMa3.1-8B (extractive)	<u>37,06</u>	29,06	32,58	<u>44,07</u>	34,55	38,74
GliNER (zero-shot)	30,6	<u>50,79</u>	<u>38,19</u>	35,33	<u>58,64</u>	<u>44,09</u>
GliNER (fine-tuned)	<b>75,15</b>	<b>63,74</b>	<b>68,98</b>	<b>82,4</b>	<b>69,9</b>	<b>75,64</b>

\* Bold and underlined represent best and second-best results respectively

- All models used in a zero-shot setting are outperformed by the fine-tuned model
- LLaMa3.1 with a generative prompt shows the worst performance
- Using an extractive prompt increases significantly the precision for the LLM

# Per-class Results

		Exact				Fuzzy			
		PER	LOC	WORK	Avg.	PER	LOC	WORK	Avg.
Precision	LLaMa3.1-8B (generative)	<u>28,97</u>	9,00	12,69	16,89	<u>29,85</u>	9,00	18,07	18,97
	LLaMa3.1-8B (extractive)	<u>56,87</u>	<u>15,85</u>	<u>15,19</u>	<u>29,30</u>	<u>61,02</u>	<u>17,07</u>	<u>28,92</u>	<u>35,67</u>
	GliNER (zero-shot)	<u>45,18</u>	<u>12,96</u>	<u>15,86</u>	<u>24,67</u>	<u>46,68</u>	<u>14,07</u>	<u>29,94</u>	<u>30,23</u>
	GliNER (fine-tuned)	<b>89,75</b>	<b>81,25</b>	<b>44,50</b>	<b>71,83</b>	<b>92,00</b>	<b>81,25</b>	<b>63,50</b>	<b>78,92</b>
Recall	LLaMa3.1-8B (generative)	59,76	16,39	<u>31,28</u>	35,81	61,59	16,39	44,55	40,84
	LLaMa3.1-8B (extractive)	36,18	21,31	14,69	24,06	38,82	22,95	27,96	29,91
	GliNER (zero-shot)	<u>60,98</u>	<u>57,38</u>	25,11	<u>47,82</u>	<u>63</u>	<u>62,29</u>	<u>47,39</u>	<u>57,56</u>
	GliNER (fine-tuned)	<b>72,97</b>	<b>63,93</b>	<b>42,18</b>	<b>59,69</b>	<b>74,80</b>	<b>63,93</b>	<b>60,19</b>	<b>66,31</b>
F1	LLaMa3.1-8B (generative)	39,02	11,63	18,06	22,9	40,21	11,63	25,72	25,85
	LLaMa3.1-8B (extractive)	44,22	18,18	14,94	25,78	47,45	19,58	28,43	31,82
	GliNER (zero-shot)	<u>51,9</u>	<u>21,15</u>	<u>19,45</u>	<u>30,83</u>	<u>53,63</u>	<u>22,96</u>	<u>36,7</u>	<u>37,76</u>
	GliNER (fine-tuned)	<b>80,49</b>	<b>71,56</b>	<b>43,3</b>	<b>65,11</b>	<b>82,51</b>	<b>71,56</b>	<b>61,8</b>	<b>71,96</b>

\* Bold and underlined represent best and second-best results respectively

# Key Insights

- Large pre-trained models like LLaMa require significant adaptation to handle the nuanced needs of NER on historical and literary texts
- All models are less prone to errors with respect to the person class
- Literary works remain the most challenging class due to varied references and complex naming conventions
- Challenges of lexical variations and abbreviations in surface forms (e.g., *Ven.* for Venezia)



# Future Work

- **IE Benchmarking:** Collect further annotated documents to benchmark NER and EL on Italian texts from different periods
- **Entity Linking:** Extend the models' capabilities to not only recognize but also link entities to Wikidata elements
- **Model Improvement:** Develop methodologies for integrating semantic, logical, and statistical approaches to refine NER and EL predictions



# Useful Links

Leopardi Knowledge Graph:

- [https://sntcristian.github.io/leopardi\\_kg/](https://sntcristian.github.io/leopardi_kg/)

Zibaldone Entity Linking Dataset:

- <https://zenodo.org/records/14103094>

Cambridge semi-diplomatic edition:

- <https://cudl.lib.cam.ac.uk/collections/leopardi/1>

DigitalZibaldone:

- <https://digitalzibaldone.net/>

Biblioteca Italiana:

- <http://www.bibliotecaitaliana.it/>

Digital Library BNN:

- <https://dl.bnnonline.it/handle/20.500.12113/4758>

Leopardi Ecdosys:

- <https://leopardi.ecdosys.org/it/Home/>

WikiLeopardi:

- [https://wikileopardi.altervista.org/wiki\\_leopardi/index.php?title=Wiki\\_Leopardi](https://wikileopardi.altervista.org/wiki_leopardi/index.php?title=Wiki_Leopardi)



# Thank you!



Cristian Santini, Laura Melosi, Emanuele Frontoni

email: [c.santini12@unimc.it](mailto:c.santini12@unimc.it)