LLMs and Symbolic Meaning Representations

Perspectives and Questions

Jan Hajič hajic@ufal.mff.cuni.cz



June 26, 2025



Outline

- Generative Large Language Models
 ... And the OpenEuroLLM project
- Meaning Representations
 - Latest developments: PDT, UMR and beyond
- How are LLMs related to Meaning Representations
 - Basic Research Questions
- Conclusions



d esentations

LLMS



What is a (Generative) Large Language Model

- Known to the public primarily as conversational LLM (e.g. ChatGPT, Llama-3.3-70B-Instruct)
- Technology
 - **Deep Neural Networks**
 - Trained from data (texts) Machine Learning \bigcirc
 - Basic function: generate next word (segment, token) based on (long) sequence of previous words (tokens)
 - In interactive systems: start with a user "prompt"
 - Can be up to a million words (in some LLM systems)



What is a (base, foundation(al)) LLM?

- Model trained on running text only
 - i.e., not interactive
 - (i.e., cannot answer questions [well])
 - Can be monolingual, multilingual, include (source) code \bigcirc
 - Can be multimodal (w/suitably encoded images, video, etc.)
- It is a basis for applications
 - Interactive (chatbot, conversational) LLM is created by
 - fine-tuning, continuous pre-training
 - human interaction annotated data, relevance rating, etc.

5



How large is a Large Language Model?

- Model size is specified as
 - Number of parameters (weights), in millions (M) or [U.S.] billions (B) \bigcirc
 - Weight is a "real" number in certain precision (from 32 down to 1.58 bit)
 - From that, byte size can be calculated: (1B weights à 8 bit = 1 GB)
- Known model sizes (open-weight models)
 - Llama 3.1: 405 B parameters (META)
 - Llama 3.3: 70 B parameters \bigcirc
 - Quantized (lower precision than original) e.g. to 6 bits: 53 GB size)
 - For inference (",runtime"): 1 or more GPU cards
 - Context size matters: takes a large proportion of GPU card's memory
- Model training:
 - Number of parameters fixed (in the standard setting) \bigcirc
 - Different data (text) sizes (in words/tokens: tokenization very important) \bigcirc
 - Llama 3.1: 15 T (trillion) tokens



OpenEuroLLM

Our goal:

Open **Multilingual** European Generative Foundational LLM

- **Open Source** (in full) ightarrowincluding fully inspectable data
- 37+ languages ullet
 - EU + associated (+ business)
- **High-quality** ullet
 - standard and native benchmarks
- Compliant with EU regulations igodol





















































Wider context

- Programme: Digital Europe (25/50% co-funding), 3 years
- Set of AI-06 calls (projects started Jan-Mar 2025):
 - Two large projects: OpenEuroLLM and LLMs4EU
 - Coordination (ALT-EDIC4EU), total <u>~80 mil. EUR + HPC</u>
 - Strong cooperation (Deploy AI, TAILOR, TrustLLM, HPLT, ...)
- Goals
 - Develop open LLMs, including conversational
 - Adapt them to applications in all areas, from commerce to e-government and education
 - Contribute to EU's digital sovereignty





Collaborations, Compute, and Data

- Open source community members
 - Experts on LLMs (incl. from non-EU ones, informal)
 - LAION, open-sci, Common Crawl, ...
 - Experts on legal issues

OPEN

EURO

- Computing power: 5 EuroHPC centers on board (project partners)
 - Now: 3m GPU hours on Leonardo (CINECA), 1.5m GPU on LUMI-G (CSC)
 - For generating synthetic data, MoE experiments, multilinguality
- Data (w/CommonCrawl, Internet Archive, OpenWebSearch)
 - From previous projects (HPLT) and other sources 37+ languages!
 - Cleaning, language ID, topic detection essential



artners) LUMI-G (CSC) nguality n)

Evaluation and Benchmarking

- For initial experiments:
 - Standard benchmarks for base models
- Project longer-term goal
 - Benchmarks for all languages in native form
 - i.e., manually translated or inspected, incl. contents
- **Continuous evaluation**
- Tests for evaluation data purity
 - I.e., not used in training/SFT/...
- Models released based on evaluation results





Meaning Representation(s)



15

Meaning Representation is...

- Symbolic system for describing "meaning"...
 - Usually in the form of a graph
 - Nodes: units of meaning
 - Edges: relation between such units
 - Attributes additional information at nodes and/or edges
- In related to the (surface, running) text that expressed that
 - Aligned to text
 - At a sentence level only
 - With word/token granularity ("anchored")



edges essed that

There are many Meaning Representations...

- Meaning representations vary along many dimensions
 - How meaning is connected to text
 - Anchoring, alignment, multi-layer vs. text-span only
 - Relationship to logical and/or executable form
 - Mapping to Lexicons/Ontologies
 - General, task-oriented

- Relationship to discourse and discourse-like phenomena
 - Including co-reference, information structure, scoping, etc.
 - ... and other inter-sentential relations



17

There are many Meaning Representations...

Μe		Alignment	Logical Scoping & Interpretation	Ontologies and Task-Specific	Discourse-Level
•	DRS (Groningen / Parallel)	Compositional /Anchored	Scoped representation (boxes)	Rich predicates (WordNet), general roles	Can handle referents, connectives
•	MRS	Compositional /Anchored	Underspecified scoped representation	Simple predicates, general roles	n/a
•	UCCA	Anchored	Not really scoped	Simple predicates, general roles	Some implicit roles
	Prague (PDT) Tectogrammatical Representation Layer	Anchored	Not really scoped with exceptions (negation)	Rich predicates, semi-lexicalized roles	Rich multi-sentence conference, discourse
	AMR	Unanchored except ZH	Not really scoped yet	Rich predicates, lexicalized roles	Rich multi-sentence coreference
•	UMR	Anchored	Scoped (quantifiers, negation)	Rich predicates, lexicalized roles	Rich multi-sentence conference, discourse



Prague Dependency Treebank (PDT) in a nutshell







"TR" (Tectogrammatical layer)

Surface dependencies (not UD but convertible)

Morphology

osuffice[-you] o. stěží vystačíte stěží-2 vystačit VB-P---2P-AAP-- 7.



V nedělních parlamentních volbách v Estonsku získal podle včerejších předběžných výsledků nejvíce hlasů blok Vlast, jehož prezidentským kandidátem byl Lennart Meri.

'In Sunday's parliamentary elections in Estonia, according to yesterday's preliminary results, the Homeland bloc, whose presidential candidate was Lennart Meri, won the most votes.'

(borrowed from the PDiT-EDA 1.0 corpus; English glosses added).



Ō



V nedělních parlamentních volbách v Estonsku získal podle včerejších předběžných výsledků nejvíce hlasů blok Vlast, jehož prezidentským kandidátem byl Lennart Meri.

'In Sunday's parliamentary elections in Estonia, according to yesterday's preliminary results, the Homeland bloc, whose presidential candidate was Lennart Meri, won the most votes.'

(borrowed from the PDiT-EDA 1.0 corpus; English glosses added).



ID

Uniform Meaning Representation (UMR)





Uniform Meaning Representation (UMR)





Uniform Meaning Representation (UMR)





Beyond Predicate-Argument Structure

- **Current Predicate-Argument lexicons**
 - PropBank (AMR/UMR), (PDT-)Vallex (PDT and similar)
- Known issues
 - Argument labels not always "semantic"
 - ?semantics of PAT in PDT, Arg2 in PropBank
 - Synonymy?
 - Is there a difference, in a particular context, among "inform", "announce" and "tell"? If yes, in what exactly?
 - Not having it makes inference more difficult
 - Hierarchy (IS-A, or general/specific relation)?
 - More general terms often used but actual meaning is more specific





SynSemClass ontology of eventive types

- Class ~ eventive concept ("to have something in possession")
 - Class members: words (senses) with argument structure and roles
 - A.k.a. synonyms
 - In multiple languages (concept is "language independent")





SynSemClass ontology of ever

- Class ~ eventive concept ('
 - Class members: words (sense
 - A.k.a. synonyms
 - In multiple languages (c

own (ev-w2176f1) vlastnit (v-w7650f1) besitzen (VALBU-ID-400394-1) poseer (AnCora-ID-poseer-1) Class ID: vec00348 Roleset: Asset controller^{def.}; Asset^{def.} Classmembers: hold (EngVallex-ID-ev-w1601f7) ACT: PAT **own** (EngVallex-ID-ev-w2176f1) ACT; PAT Own sth **possess** (EngVallex-ID-ev-w2340f1) ACT: PAT **Own sth** držet (PDT-Vallex-ID-v-w839f3) ACT; PAT patřit (PDT-Vallex-ID-v-w3411f2)

PAT; ACT Own sth

vlastnit (PDT-Vallex-ID-v-w7650f1)
ACT; PAT
Own sth

besitzen (VALBU-ID-400394-1) VA0: VA1





ssion")





SynSemClass ontology of eventive types

- Class ~ eventive concept ("to have something in possession")
 - Class members: words (senses) with argument structure and roles
 - A.k.a. synonyms
 - In multiple languages (concept is "language independent")
- **Definitions**, examples
- Links to existing lexical-semantic resources
 - FrameNet, WordNet, VerbNet, OntoNotes and similar in other languages





SynSemClass ontology of eventi

- Class ~ eventive concept ("to
 - Class members: words (senses)
 - A.k.a. synonyms
 - In multiple languages (cor
- Definitions_examples
- Links to existing lexical-sema
 - FrameNet, WordNet, VerbNet,

own (ev-w2176f1)

Class ID: vec00348^{def.}

Roleset: Asset_controller^{def.}; Asset^{def.}

Classmembers:

	hold (EngVallex-ID-ev-w1601f7					
	ACT; PAT					
	FN: NM					
	WN: hold#29; hold#4; hold#9					
	ON: hold#7					
	VN: own-100.1					
	PB: NM					
	own (EngVallex-ID-ev-w2176f1)					
	ACT; PAT					
	Own sth					
	FN: Possession/own.v					
	WN: own#1					
	ON: own#1					
	VN: own-100.1					
	PB: own/own.01					
possess (EngVallex-ID-ev-w23						
	ACT; PAT					
Own sth						
	FN: Possession/possess.v					
	WN: possess#1; possess#2					
	ON: possess#1					
	VN: own-100.1					

PB: possess/possess.01









SynSemClass ontology of eventive types

- Class ~ eventive concept ("to have something in possession")
 - Class members: words (senses) with argument structure and roles
 - A.k.a. synonyms
 - In multiple languages (concept is "language independent")
- Definitions, examples
- Links to existing lexical-semantic resources
 - FrameNet, WordNet, VerbNet, OntoNotes and similar in other languages

Hierarchy in classes

More general class / more specialized classes





SynSemClass ontology of eventive types of eventi

- Class ~ eventive concept ("to have
 - Class members: words (senses) with a
 - A.k.a. synonyms
 - In multiple languages (concept is
- Definitions, examples
- Links to existing lexical-semantic re
 - FrameNet, WordNet, VerbNet, OntoNot
 - Hierarchy in classes
 - More general class / more specialized c





	hic_0
(2 / 1 / 176)	hic_1
(4 / 1 / 154)	hic_1_1
1 / 64)	hic_1_1_4
1 / 35)	hic_1_1_4_3
_cause-and-effect (07171)	hic_1_1_4_3_4
_owner (4 / 1 / 9)	hic_1_1_4_3_1
nce or Availability (2 / 1 / 4)	hic_1_1_4_3_1_2
or Loss (0 / 1 / 1)	hic_1_1_4_3_1_3
or Source (0 / 1 / 1)	hic_1_1_4_3_1_4
ssion or Ownership (0 / 2 / 2)	hic_1_1_4_3_1_1
_purpose (0 / 1 / 1)	hic_1_1_4_3_2
_situation (5 / 2 / 23)	hic_1_1_4_3_3
19)	hic_1_1_4_1
1 / 9)	hic_1_1_4_2
ianent States (2 / 0 / 7)	hic_1_1_1
permanent States (3 / 0 / 80)	hic_1_1_3
ent States (0 / 2 / 2)	hic_1_1_2
/ 21)	hic_1_2
	hic_2
	hic_2_9
	hic_2_6
	hic_2_4
0 / 255)	hic_2_2
	hic_2_7
/ 12)	hic_2_3
' 0 / 319)	hic_2_1
/ 298)	hic_2_5
	hic_3
))	hic_3_3
	hic_3_2
	hic_3_1
272)	hic_4

LLMs vs. (and?) Meaning Representations



(Computational) Linguistics, NLP and LLMs

- Is NLP "solved" (by LLMs)?
 - No, not really, not yet but there is great progress
 - Still hard problems remain
 - reasoning, explainability, interpretability; low-resourced languages
 - Can Semantic Representations be used with(in) LLMs to improve NLP?
- Is (Computational) Linguistics "solved" (by LLMs)?
 - Not at all
 - We can ask LLMs questions about language
 - Answers come from the texts they were trained on ...
 - In their "introspection"
 - Hard questions unresolved



nguages ve NLP?

Open questions

- Linguistics
 - Language structure (is there any...?)
 - What are symbolic representations actually telling us?
 - About morphology, syntax, meaning, ...
- (Language) Learning
 - How do we learn language (mother tongue, 2nd)? Anything LLMs can teach us?
- Relation between language and the world around us
 - How is our knowledge (memory) structured? Any parallels with LLMs?
 - Relation between perception of language, vision and other senses
 - Why are we describing language as graphs or formulas?
 - What exactly is "grounding" (in perception, communication, society)?



LLMs vs. computational models

- How do they represent the language they were trained on?
- How do they generalize (in the sense humans do)? Do they?
- How do they learn about "concepts", represent and reason over them?
- Can we learn anything from comparing humans and LLMs?
 - How can we make such comparisons?
 - Brain level, symbolic level, using logic, philosophy, ...?
 - Is it fair comparison? (different learning mechanisms, language x other modalities, ...)
- Should we collaborate with psycholinguists, neurolinguistics, cognitive scientists, logicians, philosophers? [yes, of course; but how exactly?]
- Where LLM biases come from and how they differ from ours?



d on? o they? ason over them? LMs?

x other modalities, ...) uistics, cognitive t how exactly?]





Questions?



https://ufal.mff.cuni.cz/node/2540 UFAL MFF UK – UMR project





Supported by the project OpenEuroLLM, GA No. 101195233, ALT-EDIC4EU, GA No. 101195344, Digital Europe Programme by *European Commission and co-funded by the JU subprogramme of* the MEYS CR and other MEYS CR and CSF programs.



Co-funded by the European Union



https://ufal.mff.cuni.cz/grants/lusyd UFAL MFF UK – LUSyD project



Open Source and Community

- Open Strategic Partnership Board (Strategic advisory role)
 - Open source community members
 - Experts on LLMs (incl. from non-EU ones)
 - Former commercial and/or open source model developers
- Experts on legal issues
- Informal cooperations
 - Data side: CommonCrawl, Internet Archive EU, OpenWebSearch
 - Open source models community
 - EuroLLM (Univ. of Edinburgh UK, UnBabel Portugal)
 - LAION, open-sci, ...





Computing facilities

- 5 EuroHPC centers on board (project partners)
 - Technical expertise jumps start using the respective facilities
- Some compute available from previous projects
- Participation in EuroHPC calls in 2025
 - In line with project plan for the rest of 2025
- ", Strategic" allocations in the future ("STEP")
 - Using current facilities & new in AI Factories (2026/2027)
 - Just received 3m GPU hours for May-Nov. 2025 on Leonardo (CINECA)
 - For generating synthetic data
 - ... and 1.5m GPU hours at CSC, on LUMI-G
 - Testing data staging, multilingual training, MoE (~scaling laws)





38

Data for 37+ languages

- Using available Open Source data
 - HPLT 2.0 (HPLT 3.0, July-Aug 25), Fineweb2, Cultura-X, ...
 - Mixtures to be experimentally determined
 - Ultimate (re)sources: CommonCrawl, Internet Archive, IA Europe
 - OpenWebSearch negotiations ongoing
- Focus on low-resource languages for additional data
 - Incl. specific cases for very similar languages
- Additional data for
 - Fine-tuning, instruction-tuning, reasoning
 - ... if necessary for benchmarking





