

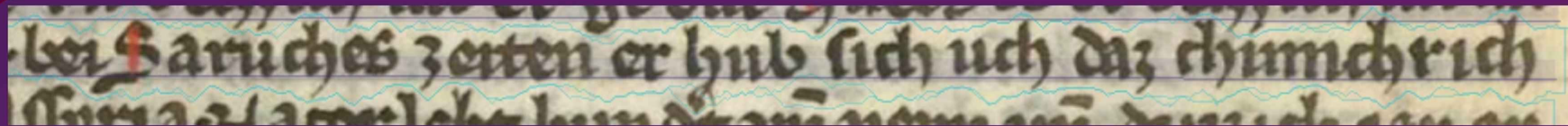
# HYDRA.



## (UNIVERSAL) TAGGER-LEMMATISER WITH DEEP LEARNING AND PARALLEL COMPUTING

### PREMODERN SPELLING

IN THE FIELD OF HISTORICAL COMPUTER LINGUISTICS WE INVARIABLY FACE A DIFFICULTY RELATED TO RELIABLE (PRE-)PROCESSING, MORE PRECISELY SPEAKING, TAGGING AND LEMMATISATION OF PREMODERN TEXTS. THE PROBLEM IS THE NATURE OF PREMODERN SPELLING, BOTH IN LATIN AND IN THE VERNACULAR LANGUAGES OF EUROPE, EXPOSING A HIGH DEGREE OF VARIANCE (BETWEEN REGIONS AND EVEN BETWEEN PARTICULAR SCRIBES, AS FAR AS SCRIBAL CULTURE IS CONSIDERED) THAT INEVITABLY IMPEDES THE AUTOMATIC TEXT PROCESSING ON SCALE. NOT ONLY A SINGLE WORD MAY HAVE BEEN SPELLED DIFFERENTLY, BUT THE BOUNDARIES BETWEEN WORDS AND MORPHEMES WERE CHANGING FROM ONE MANUSCRIPT TO ANOTHER. THAT FLUIDITY IN ORTHOGRAPHIC PRACTICE PROBABLY FORMS THE BIGGEST IMPEDIMENT TO NATURAL LANGUAGE PROCESSING OF PREMODERN HISTORICAL TEXTS.

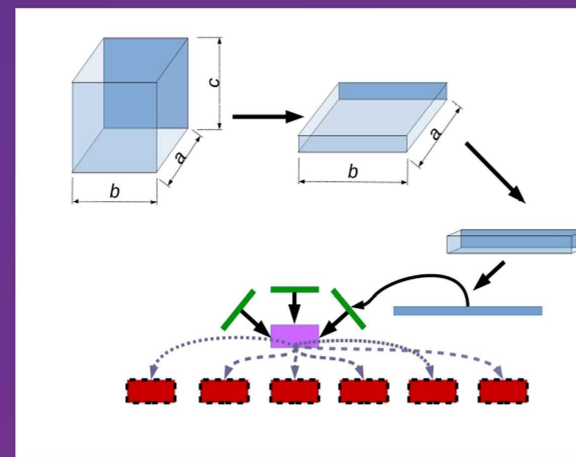
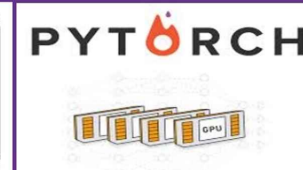


**MULTI-TOKEN SINGLE LABEL PROBLEM.** IN THIS MIDDLE HIGH GERMAN SENTENCE *BEI SARUCHES ZEITEN ER HUB SICH UCH DAS CHUNICHRIK* [IN ASSYRIA] [‘IN THE TIME OF SERUG THERE AROSE ALSO THE KINGDOM [OF ASSYRIA]’] THE REFLECTIVE VERB *ER HUB SICH* [AROSE, EMERGED] IS PRESENTED BY TWO TOKENS IN THE DIPLOMATIC TRANSCRIPTION WIDELY USED IN MEDIEVAL STUDIES.



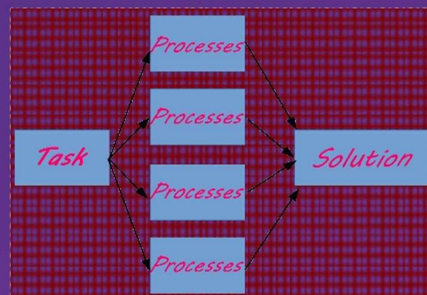
**MULTI-LABEL SINGLE TOKEN PROBLEM.** IN THIS MIDDLE LOW GERMAN SENTENCE *DAT HAUESTU GHEDAN* [‘YOU DID IT/HAVE DONE IT’] THE TOKEN *HAUESTU* IS ACTUALLY THE ENCLITIC FORM OF *HAUEST TU* AND NEEDS TWO LABELS (PART-OF-SPEECH TAGS AND LEMMAS FOR BOTH THE VERB *HAUEST* ‘(YOU) HAVE’ AND THE PERSONAL PRONOUN *TU* ‘YOU’).

### DEEP LEARNING

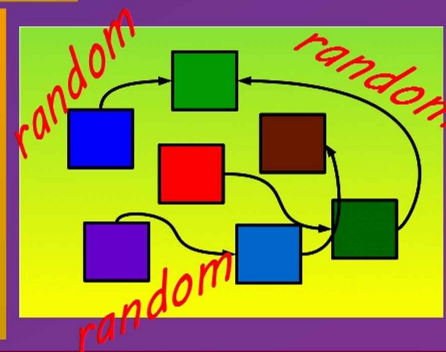


MULTI-LAYER ARTIFICIAL NEURAL NETWORKS AKA DEEP LEARNING IS A FIELD OF RESEARCH AIMING AT COMPLEX LEARNING METHODS THAT CAN EFFICIENTLY ANALYSE HUGE AMOUNT OF DATA AND SOLVE COMPLICATED NON-TRIVIAL TASKS. BY REPRESENTING WORDS AND WHOLE PHRASES OR SENTENCES AS EMBEDDING TENSORS WE COMPRESS AND MINE RELEVANT LINGUISTIC INFORMATION CONCEALED, FOR INSTANCE, IN AFFIXES OR WORDS COMBINATIONS. IN THIS WAY WE ENSURE HIGH ACCURACY OF A LABEL CLASSIFICATION DEPENDING BOTH ON THE WORD STRUCTURE AND THE WORD CONTEXT.

### PARALLEL COMPUTING AND GOSSIP ALGORITHM



PARALLEL COMPUTING IS AN UMBRELLA TERM FOR DIFFERENT METHODS OF SHARING DATA AND COMPUTATIONS AMONG MANY SEPARATE UNITS (PROCESSORS, GRAPHIC CARDS, MACHINES ETC.). ONE OF THE STRATEGIES FOR THE COORDINATION OF PARALLEL PROCESSES IS THE SO-CALLED GOSSIP ALGORITHM WHICH ALLOWS PROCESSES TO EXCHANGE RELEVANT INFORMATION BETWEEN EACH OTHER IN A RANDOM WAY. THE FINAL SOLUTION IS AN AVERAGE OF ALL THE PROCESSES’ SOLUTIONS THAT ARE SUPPOSED TO HAVE CONVERGED IN THE MEANTIME.



### SELECTED LITERATURE

- BLOT M., PICARD D., CORD M., THOME N. ‘GOSSIP TRAINING FOR DEEP LEARNING’. NIPS 2016 WORKSHOP, BARCELONA, SPAIN, DECEMBER 2016. [HTTPS://ARXIV.ORG/ABS/1611.09726](https://arxiv.org/abs/1611.09726)
- KESTEMONT, M., DE PAUW, G., VAN NIE, R. & DAELEMANS, W., ‘LEMMATISATION FOR VARIATION-RICH LANGUAGES USING DEEP LEARNING’. *DSH – DIGITAL SCHOLARSHIP IN THE HUMANITIES* 32:4 (2017), 797-815. DOI: [HTTPS://DOI.ORG/10.1093/LLC/FQW034](https://doi.org/10.1093/LLC/FQW034)
- KESTEMONT, M. & J. DE GUSSEM, ‘INTEGRATED SEQUENCE TAGGING FOR MEDIEVAL LATIN USING DEEP REPRESENTATION LEARNING’, *JOURNAL OF DATA MINING & DIGITAL HUMANITIES* (2017), PP. 17. SPECIAL ISSUE ON COMPUTER-AIDED PROCESSING OF INTERTEXTUALITY IN ANCIENT LANGUAGES, ED. M. BUECHLER AND L. MELLERIN.
- PIOTROWSKI M., *NATURAL LANGUAGE PROCESSING FOR HISTORICAL TEXTS*, SAN RAFAEL: MORGAN & CLAYPOOL PUBLISHERS, 2012.
- ZHANG Y., SHEN D., WANG G., GAN Z., HENAO R., CARIN L. ‘DECONVOLUTIONAL PARAGRAPH REPRESENTATION LEARNING’, CONFERENCE ON NEURAL INFORMATION PROCESSING SYSTEMS 2017, LONG BEACH. ARXIV:1708.04729

[HTTPS://GITHUB.COM/LUKASZ-G/HYDRA](https://github.com/LUKASZ-G/HYDRA)

### LEMMA BUILDING

SINCE LARGE CORPORA MAY ENCOMPASS TENS OF THOUSANDS OF LEMMAS, IT CAN BE ALSO QUITE DIFFICULT TO COMPUTE MODEL PROBABILITIES FOR ALL OF THEM, ESPECIALLY FOR THOSE LESS FREQUENT. AN EXPERIMENTAL APPROACH TO CONSTRUCT A LEMMA FROM SCRATCH IS COMPUTATIONALLY EASIER AND IS THOUGHT TO PREDICT LEMMAS ABSENT IN THE TRAINING CORPUS BY LOOKING AT ANALOGIES ALREADY OBSERVED.

