# Using deep learning paradigm
# to improve syllabic versification:
# A first approach

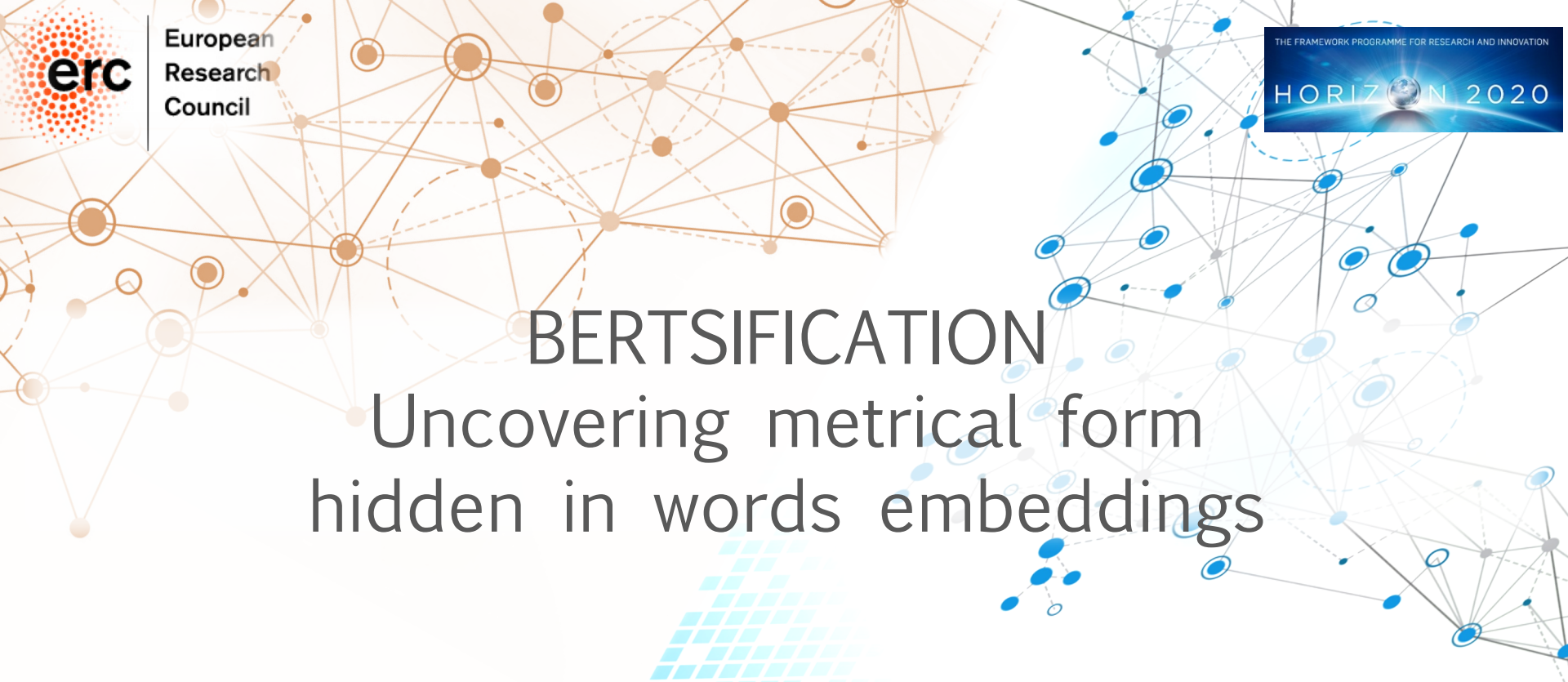Javier de la Rosa (versae@linhd.uned.es, @versae)

Postdoctoral Fellow ERC Poetry Standardization and Linked Open Data (POSTDATA) Project

Salvador Ros(sros@ssc.uned.es)

Technical Director and Co-PI ERC Poetry Standardization and Linked

Laboratorio de Innovación en Humanidades Digitales, UNED, Spain

# BERTSIFICATION
# Uncovering metrical form hidden in words embeddings

Javier de la Rosa (versae@linhd.uned.es, @versae)

Postdoctoral Fellow ERC Poetry Standardization and Linked Open Data (POSTDATA) Project

Salvador Ros(sros@ssc.uned.es)

Technical Director and Co-PI ERC Poetry Standardization and Linked

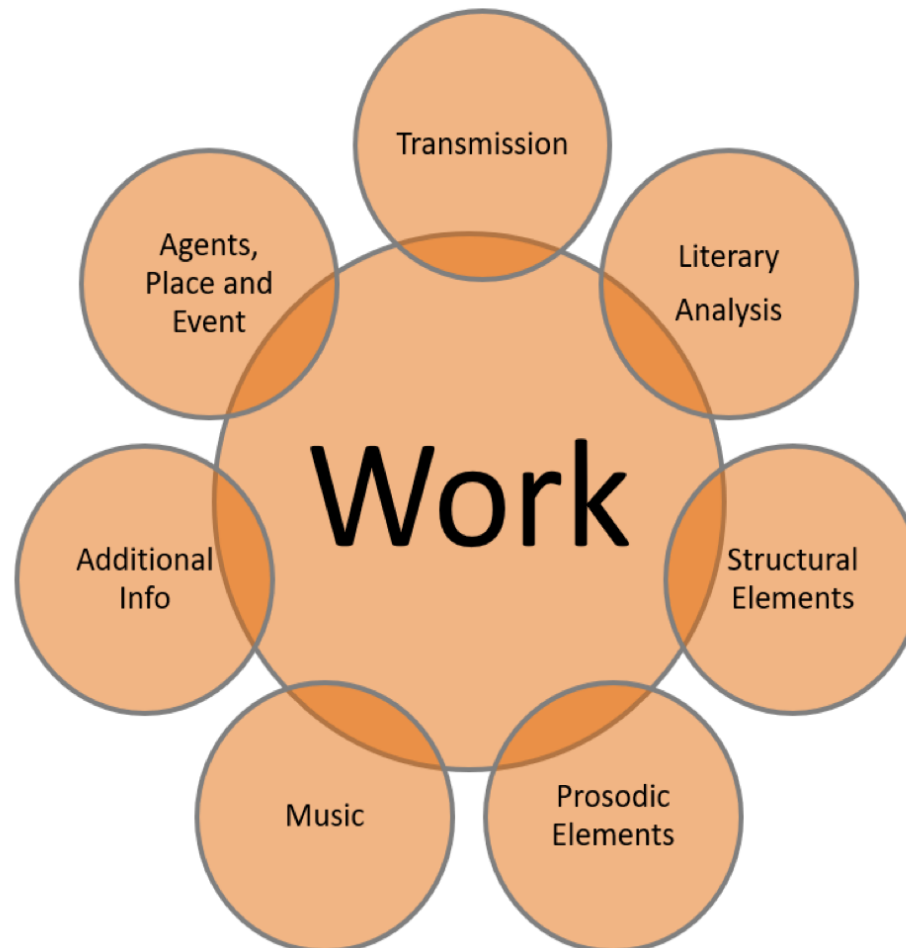Laboratorio de Innovación en Humanidades Digitales, UNED, Spain

THE FRAMEWORK PROGRAMME FOR RESEARCH AND INNOVATION

HORIZON 2020

UNED

LiNHD
LABORATORIO DE INNOVACIÓN
EN HUMANIDADES DiGiTALES

POSTDATA
Poetry Standardization
and Linked Open Data

- ERC Starting Grant
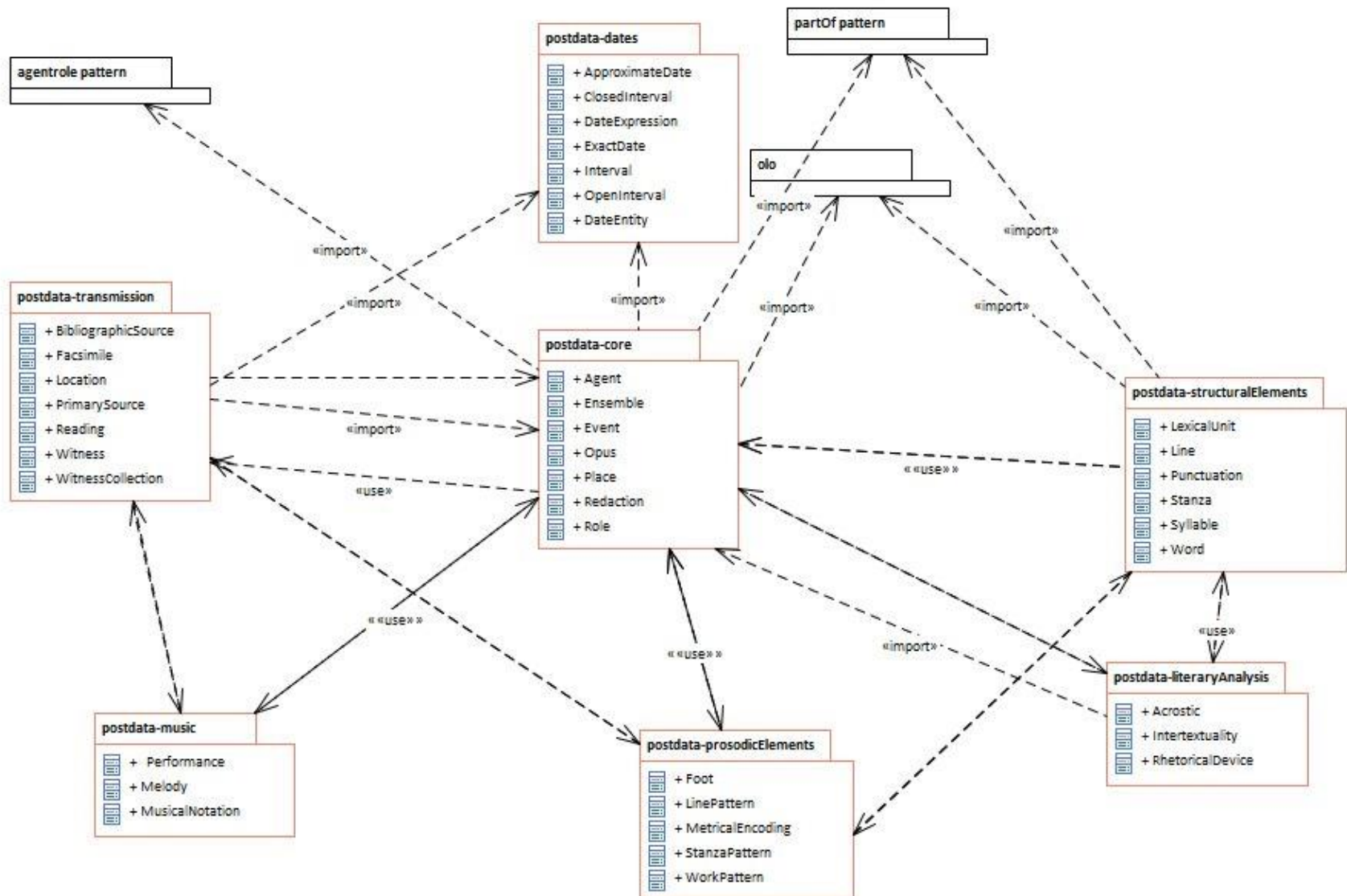
- Poetry Standardization and Linked Open Data

# POSTDATA



POSTDATA
Poetry Standardization
and Linked Open Data

- Very rich model → Hard to fill

- How can we help in filling it?

- Very rich model → Hard to fill

- How can we help in filling it?
  – Automated scansion

- Rule-based

- Inference-based

  – Generative

  – Neural scansion

  – Machine learning

# Rule-based scansion

- Syllabification


- Prosody and stress


- Patterns and exceptions

# Rule-based scansion

«Vino y ahogó sus penas»

- – Syllabification
  - Vi-no y a-ho-gó sus pe-nas

- – Prosody and stress
  - **Vi**-no‿**y**‿a‿ho-**gó**-sus-**pe**-nas

- – Patterns and exceptions
  - **Vi**-no-y‿a‿ho-**gó**-sus-**pe**-nas

# Rule-based scansion

- English


- Spanish

# Rule-based scansion for English

- Scandroid

- Zeuscansion

- Poesy (aka LitLab-poetry)


- Others: Calliope, AnalysePoems, etc.

# Rule-based scansion for English

POSTDATA
Poetry Standardization
and Linked Open Data

- Scandroid [1996, 2005]

  – iambic or anapestic meter

- Zeuscansion [2016]

  – OOV and stress-guesser

- Poesy (aka LitLab-poetry) [2018]

  – Diificult coding of new rhythmic patterns

# Rule-based scansion for Spanish

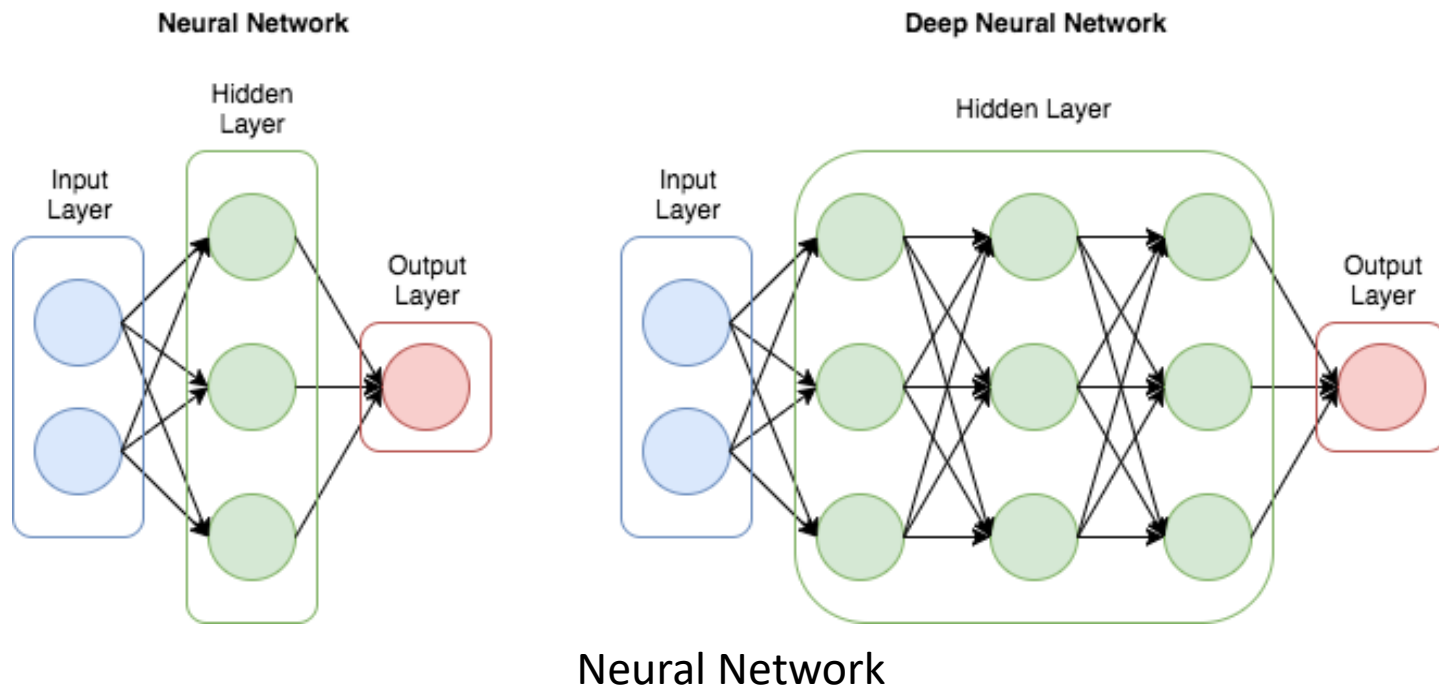- Gervás [2000]

- ADSO Scansion System [2017]

- SKAS [2017]

- Gervás [2000]

  – Logic programming

- ADSO Scansion System [2017]

  – Special focus on synaloephas

- SKAS [2017] → PoetryLab Ran·tan·plan

  – Industrial strong NLP

  – SpaCy models and API

# Neural scansion

- IXA (Spain) [2018]
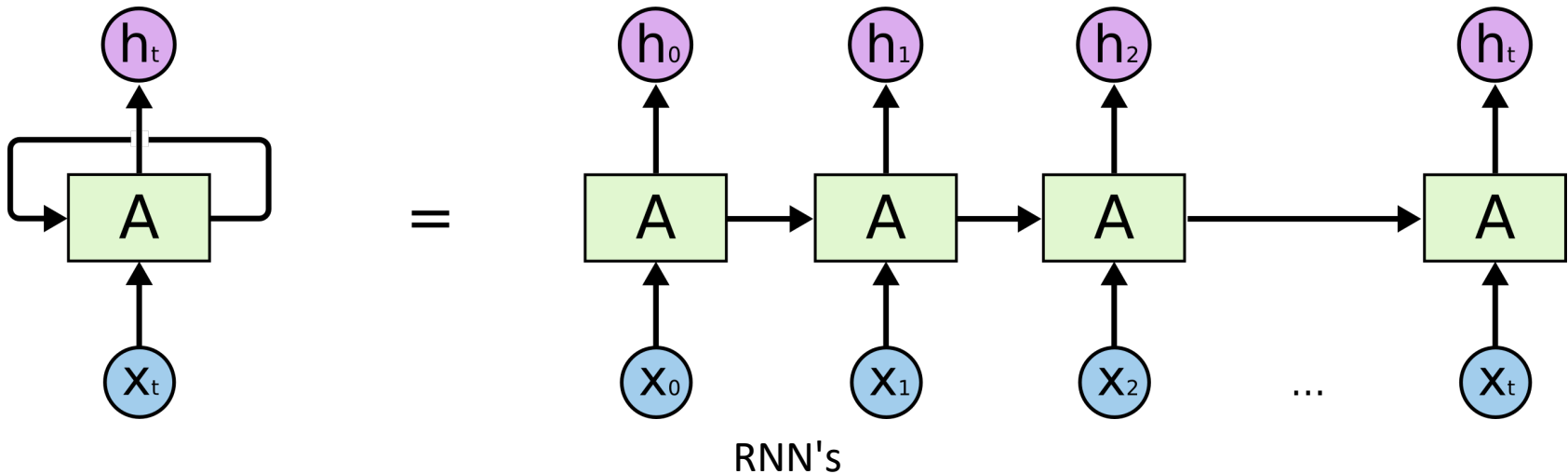  - English and Spanish

- KAIST (Korea) [2019]
  - English

# Neural scansion

- Same approach BiLSTM-{CNN,CRF}



Neural Network
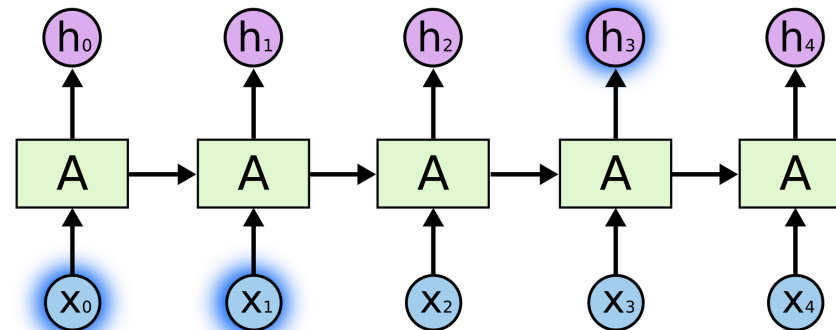
# Neural scansion

- Same approach BiLSTM-{CNN,CRF}
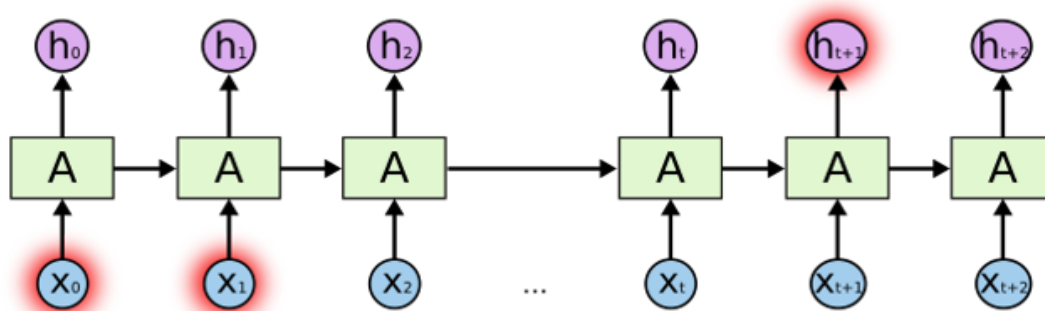


RNN's

# Neural scansion

- Same approach BiLSTM-{CNN,CRF}
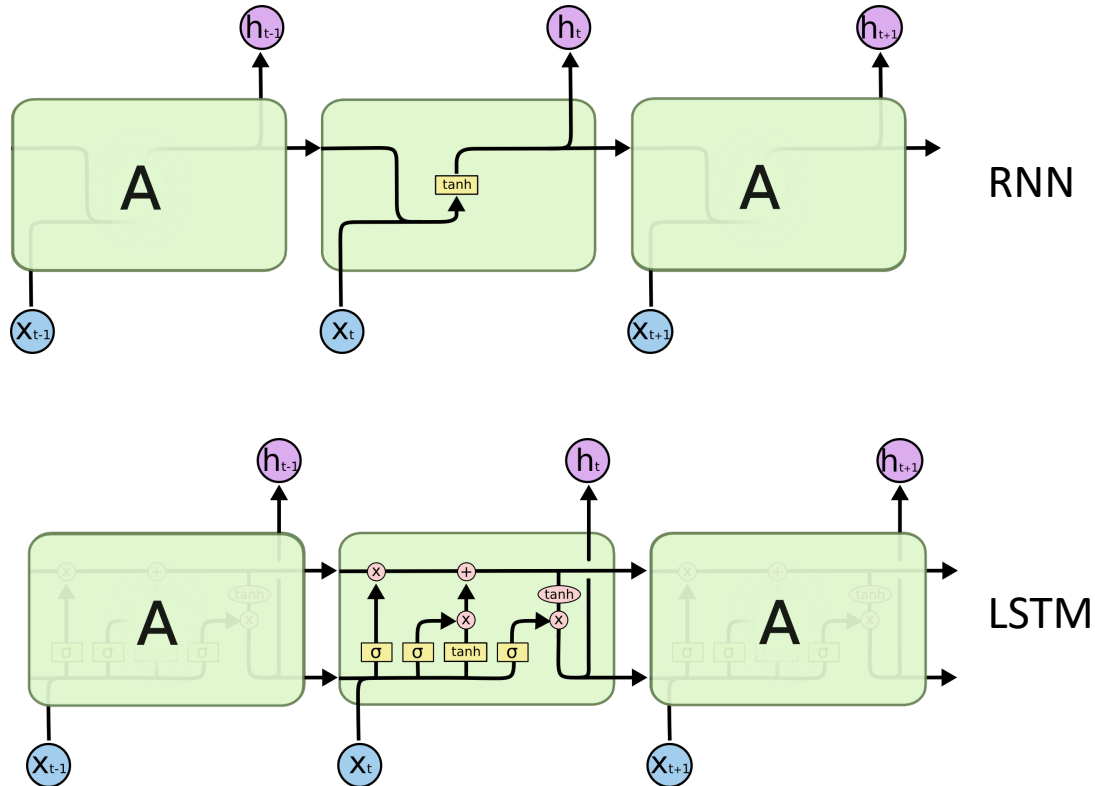


"the clouds are in the *sky*"

# Neural scansion

- Same approach BiLSTM-{CNN,CRF}



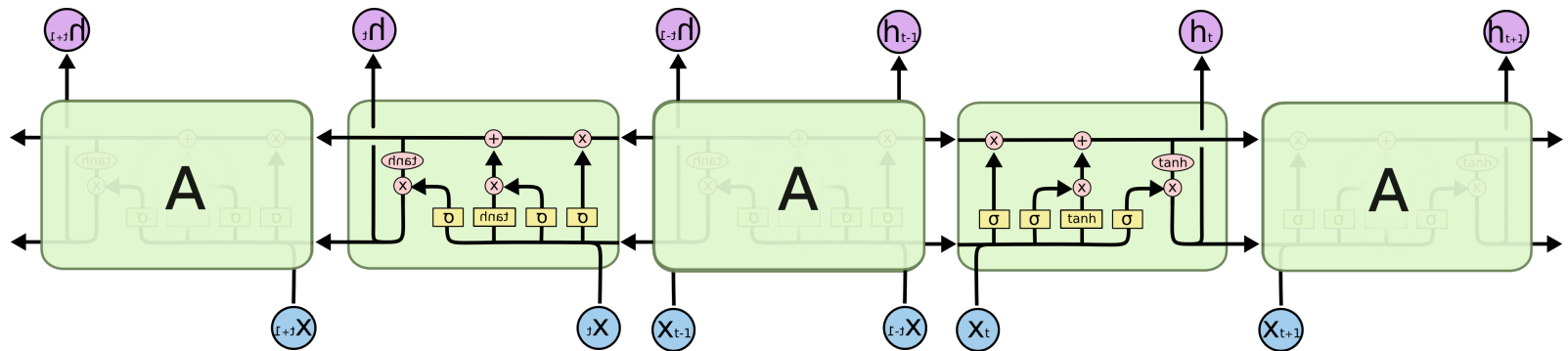"I grew up in France... I speak fluent *French*."

# Neural scansion

- Same approach BiLSTM-{CNN,CRF}



RNN

LSTM
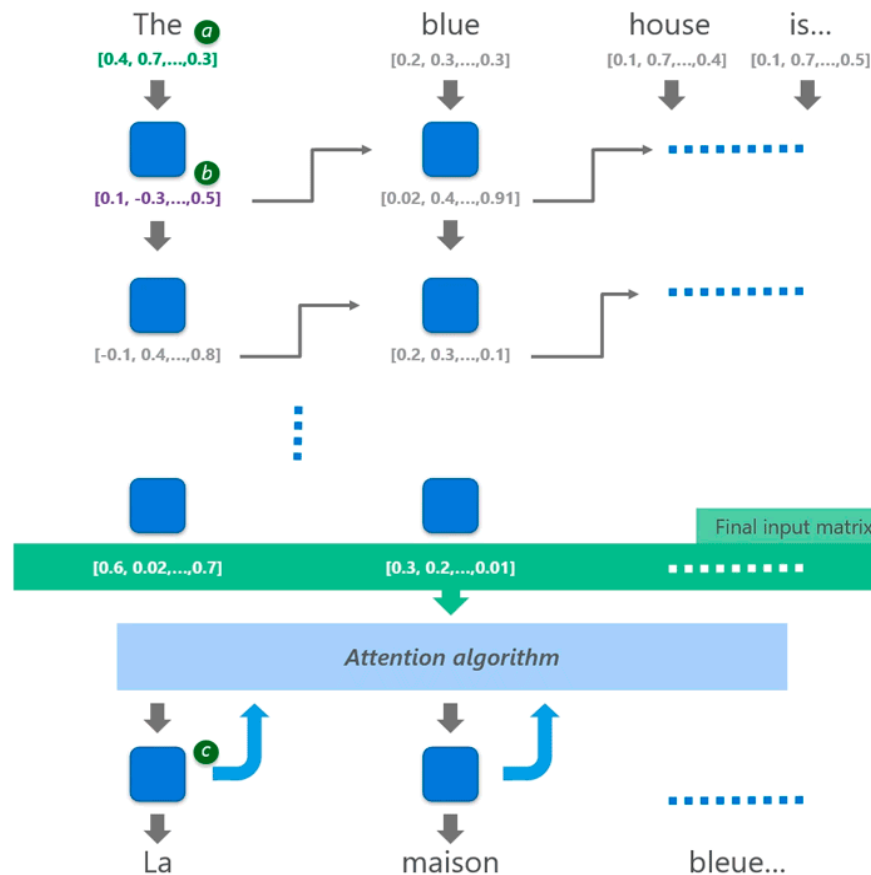
# Neural scansion

- Same approach BiLSTM-{CNN,CRF}



Bidirectional LSTM

# Neural scansion

- Attention mechanism (counts + context)

# Neural scansion

- BiLSTM with Attention

    – Our results

    – The problem of the corpus

- BiLSTM with Attention

    – Our results → 92.43% per verse (90.84% SOTA)

    – The problem of the corpus

        • Borja Navarro's annotated corpus

        • Hendecasyllables with stress in penultimate position

        • Mixed manual and automated (rule-based) annotation

POSTDATA
Poetry Standardization
and Linked Open Data

- BiLSTM with Attention

  – Our results → 92.43% per verse (90.84% SOTA)

  – The problem of the corpus

    • Borja Navarro's annotated corpus

    • Hendecasyllables with stress in penultimate position

    • Mixed manual and automated (rule-based) annotation
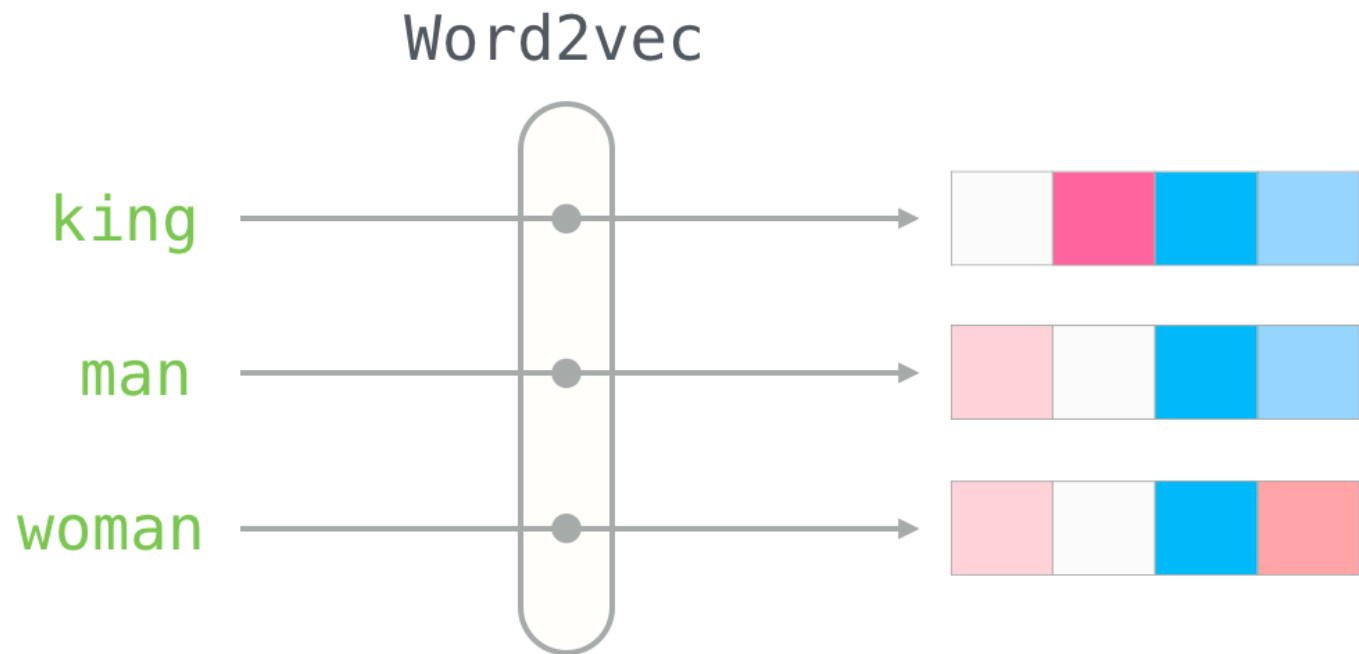
**END OF A FIRST APPROACH**

- Word embeddings

  *"You shall know a word by the company it keeps"*

  *-- J.R. Firth*

- Other applications of Attention

  - Word and sentence embeddings

  - Contextualized word embeddings

    - GPT-2 (Transformer)
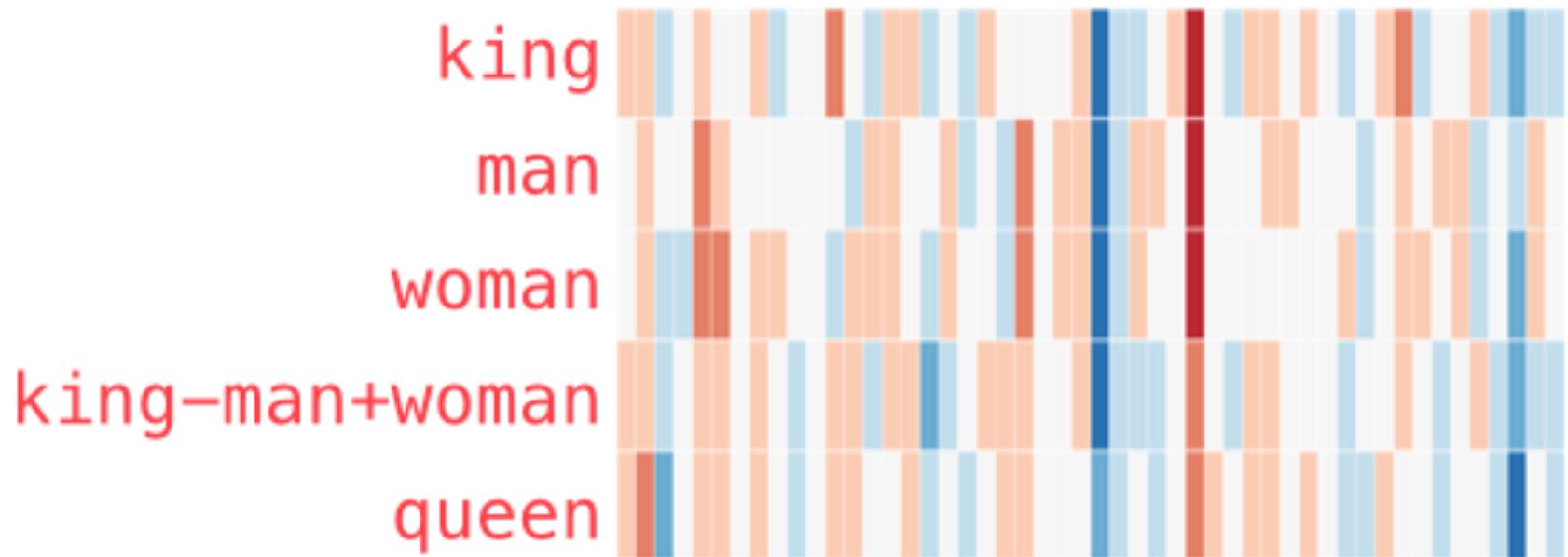
    - ULMFit

    - ELMO

    - BERT

# Neural scansion

- Word embeddings



Word2vec

king

man

woman

# Neural scansion

- Word embeddings

king − man + woman ~= queen

- Word embeddings

  "What you can cram into a single $&!#* vector: Probing sentence embeddings for linguistic properties" Alexis Conneau *et. Al.*

  – SentLen, WC, TreeDepth, TopConst, BShift, Tense, SubjNum, ObjNum, SOMO, CoordInv

# Neural scansion

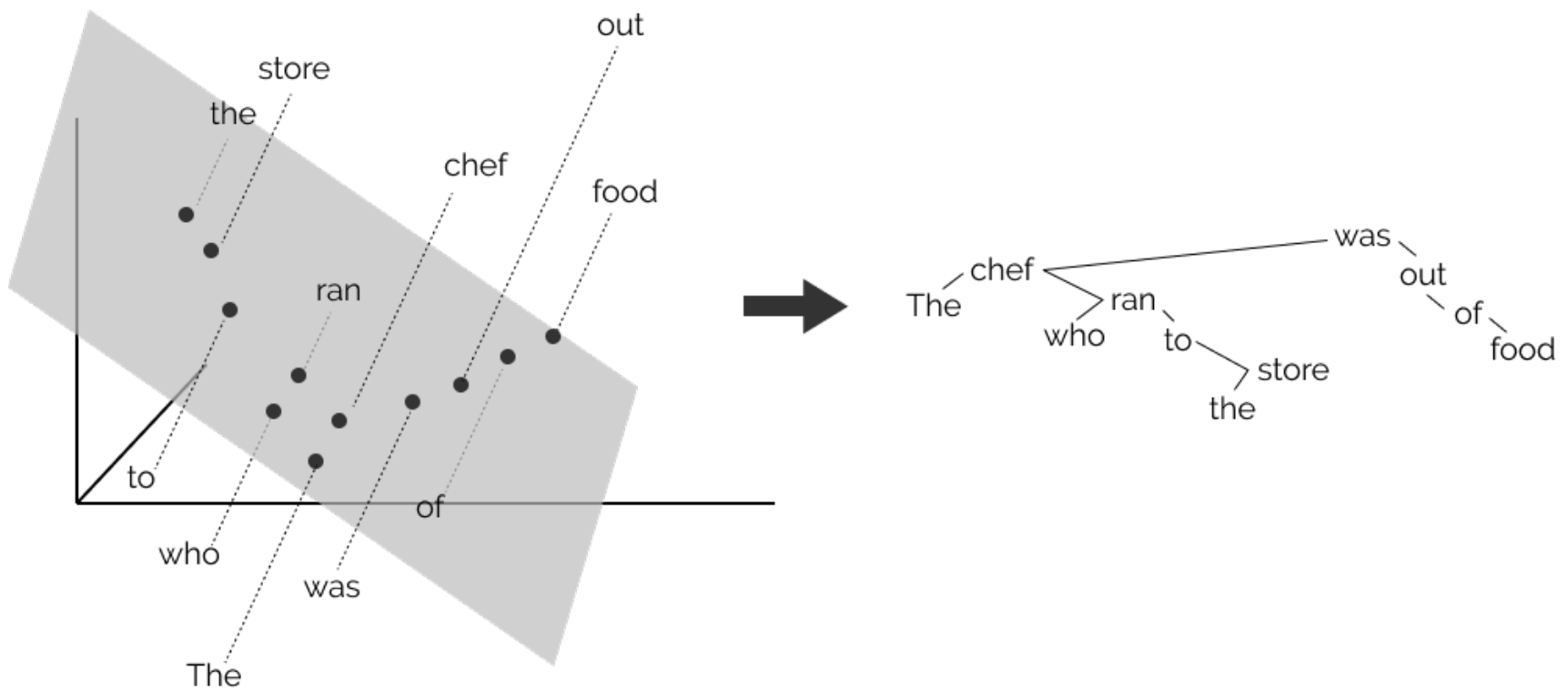- Probing contextualized sentence embeddings

# Neural scansion

POSTDATA
Poetry Standardization
and Linked Open Data

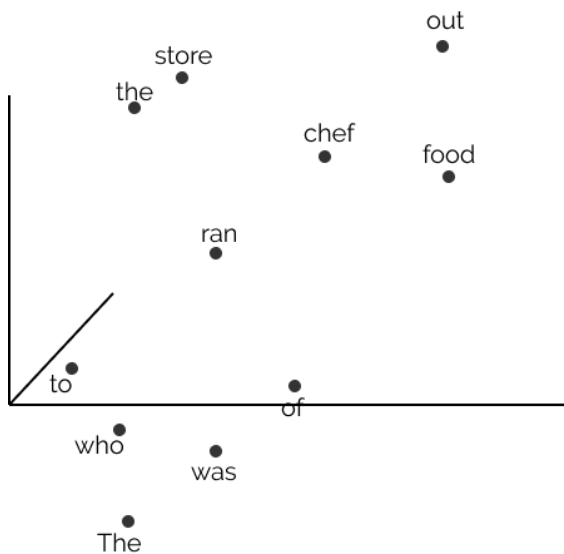- Probing contextualized sentence embeddings



Source: https://nlp.stanford.edu/~johnhew/structural-probe.html

# Neural scansion

- Other applications of Attention

  - Word and sentence embeddings

  - Contextualized word embeddings

    - GPT-2 (Transformer)

    - ULMFit

    - ELMO

    - **BERT (Bidirectional Encoder Representations from Transformers)**

- BERT

  - Multilingual and cased

  - Pre-trained model as features / Pre-train and fine-tuning

  - Successful probing

- Probing BERT (Flair) for meter information

  – **Vi**-no-y‿a‿ho-**gó**-sus-**pe**-nas

       +  -         -          +      -

       +    -          7

POSTDATA
Poetry Standardization
and Linked Open Data

- Advantages

  – Easily multi-lingual

  – Access to pre-trained models

  – No prior fixed length requirement

  – Decent **accuracy (grouped for now!)**

- Decent **enough**?

  – What does this even mean?

  – For what purpose?

  – How uncertainty should be registered or catalogued?

- Disadvantages

  – Complex construction

  – Need for pre-trained models built using massive computing infrastructures

# Conclusions

- Contextualized sentence/verse embeddings might generalize well

- Paving the way for an unsupervised multi-language scansion system

- Need better corpora!

# Questions and Thanks!

**Javier de la Rosa**
versae@linhd.uned.es
**& LINHD Team**

**http://postdata.linhd.uned.es/**
http://linhd.uned.es
@linhduned