# LDRD

**"Utilizing artificial intelligence/machine learning and automation systems to inform and optimize chemical separations"**

PI: Jasmine Hatcher-Lamarre (C-AD/MIRP)

In a nutshell:

- Designing new chemical separations takes months to years; high-throughput experimentation currently out of reach

- Literature contains useful properties and coordination trends of transition metals, lanthanides; so
    1. Extract these properties from literature with NLP, and use ML to generate new potential separation paths
    2. Use automation to efficiently implement and test those separations at scale
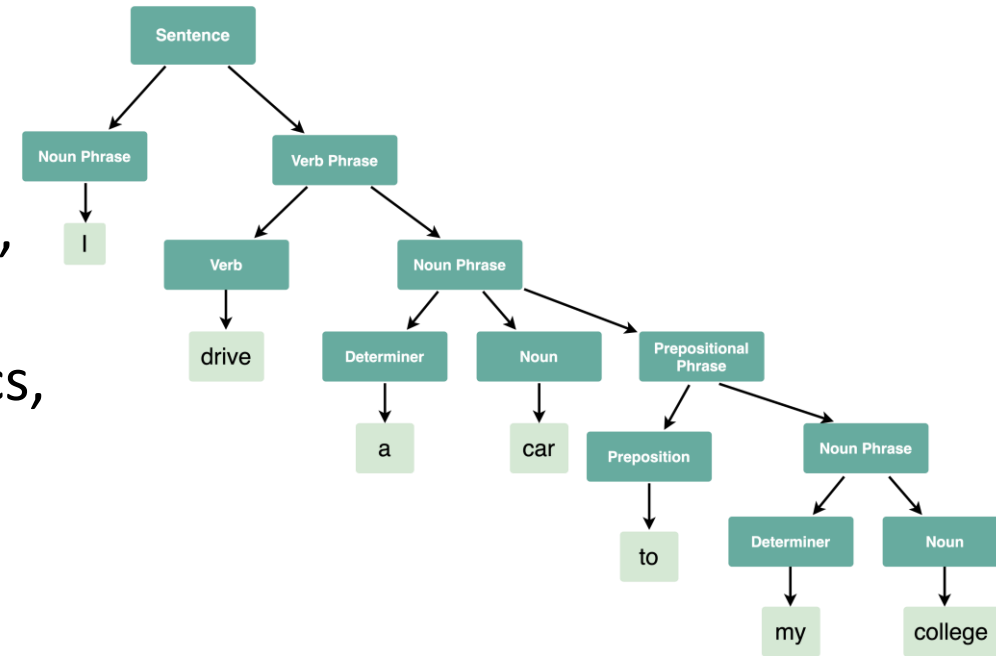
Brookhaven National Laboratory

# Extracting knowledge from literature

- Separations details in papers
  - Metals, extractants, contact/mixing time, flow rates, temperature, pH, $K_D$ constant, choice of acid, functional groups, equilibrium constants, yield, etc.
- Scaling up with NLP & ML
  - Extract & organize relevant information in paper (canonicalization)
  - Compile database of separations across many papers
  - Look for emergent patterns; synthesize new extractions
- Data relevance and variety
  - Most papers won't list all details, may cite others
  - Follow references, and broaden scope of chemical separations papers (across transition metals, lanthanides, actinides) → Trends better emerge at scale
- Building ML methods with experts in the loop
  - Include paper links and text snippets (as context) along with NLP extractions
  - Validate extractions, canonicalizations with human feedbacks
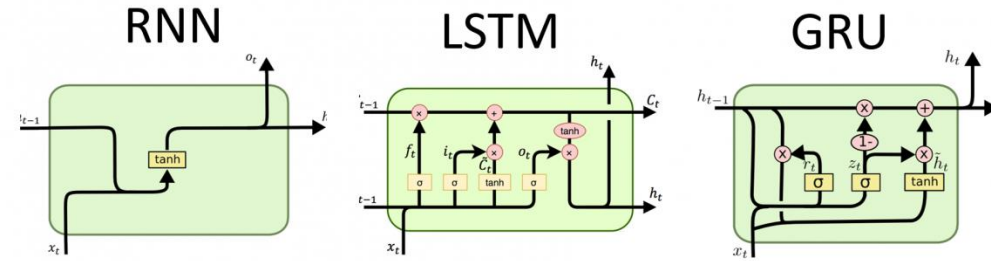  - Iteratively improve NLP and ML models

A 100 µg Th(IV) or 100 µg U(VI) was taken in 10 mL solution containing various concentrations of sulphuric acid ranging from 0.2 to 3.0 mol L$^{-1}$. The extractions were carried out in 125 mL of separating funnel. The aqueous acid solution containing metal ion was taken in 10 mL of 0.5–4.0% of 4-methyl-N–n-octylaniline in xylene for 4 min to separate the organic and aqueous phases. The metal loaded organic phase was stripped with strippant by wrist shaking of the funnel for 5 min. The 10 mL of 0.1 mol L$^{-1}$ nitric acid was used as strippant for Th(IV) while 2×10 mL acetate buffer of pH 4.5 was used as strippant in case of U(VI). The pH was adjusted with sodium hydroxide and acetic acid as per their requirement for the spectrophotometric determination. For quantification of Th(IV) xylenol orange was used as chromogenic agent while for quantification of U(VI) bromopyrogallol red was used [39].
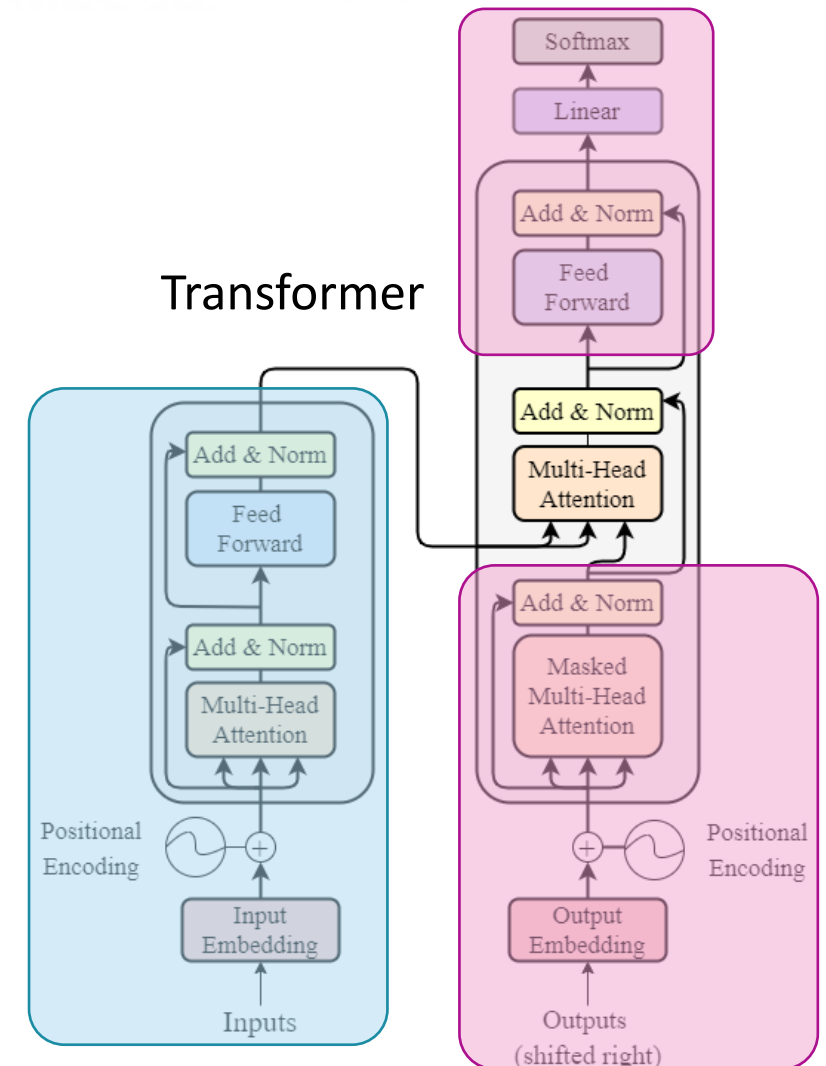
# Natural Language Processing

- Classical NLP (based on rules of linguistics)
  - Basics: parsing parts-of-speech, dependencies, syntax trees, dictionaries, grammars, etc.
  - Build up with logic rules: semantics, pragmatics, natural language understanding/synthesis
  - In practice: very brittle, poor quality
- Machine Learning for NLP
  - <u>Learn</u> language features in data-driven manner, at vast scale
  - Embeddings, language modeling, inference and generation tasks
  - New focus: Data & Models

**Brookhaven** National Laboratory
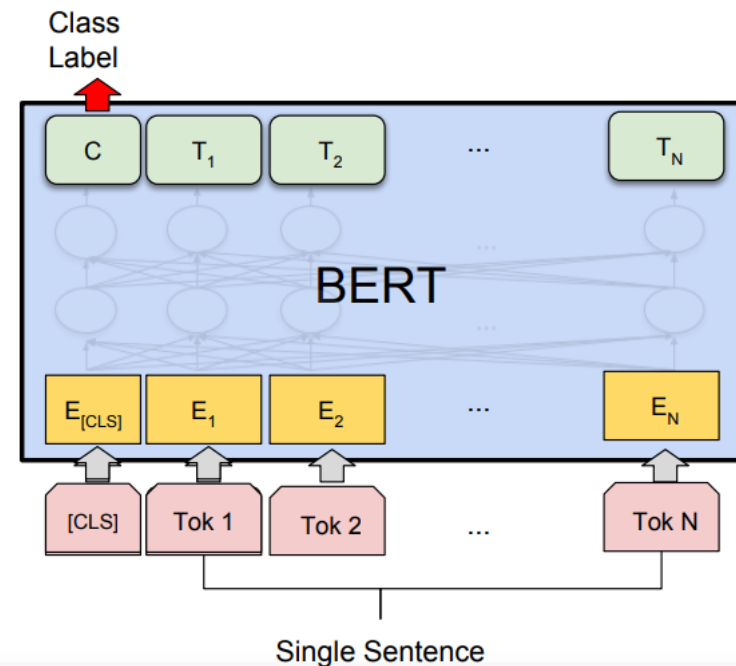
# ML Models for NLP



RNN  LSTM  GRU

- Pre-2017: Recurrent NNs (LSTMs, GRUs)

- Post-2017: Transformers[1]
  - capture long-term dependencies, scale up very well, can be pre-trained

- Encoder-type models: e.g. BERT[2]
  - Useful for building "head" onto pre-trained "body"

- Decoder-type models: e.g. GPT[3]
  - Useful for generating novel text, querying

Transformer

1. A. Vaswani, et al. "Attention is all you need." 2017.
2. J. Devlin, et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." 2019.
3. Radford, Alec, et al. "Improving language understanding by generative pre-training." 2018.

Brookhaven
National Laboratory

# Two NLP approaches

- Bottom up
  - Extract features with pre-trained encoder models (BERT[1]-like)
  - Build specialized models to recognize key elements
    - Named Entity Recognition, Text classification

- Top-down
  - Language generation with pre-trained decoder model (e.g. Galactica[2])
  - LLMs exhibit unaligned behavior
  - Prompt engineering and tuning
    - Leverage learned knowledge
    - Respond based on context from literature

J. Devlin, et al. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." 2019.
R. Taylor, et al. "Galactica: A large language model for science." 2022.

**A.8.6   Example Literature Survey Written by Galactica**

**Self-Supervised Learning, A Survey**

**Abstract:** In this paper we review self-supervised learning, a method of learning features from unlabelled data. We discuss the historical motivation behind self-supervised learning and review some of the current state of the art methods. We also look at how self-supervised learning has been used to solve a wide range of problems, including image classification, object detection, video analysis and robotics. Finally we discuss some of the challenges and future research directions for self-supervised learning.

**1 Introduction**

Deep learning has been very successful at solving many problems in machine learning, however these methods are reliant on large amounts of labelled data. In many real world applications, it is very expensive or impossible to obtain large amounts of labelled data. Self-supervised learning has emerged as a way of overcoming this issue.

[START_AMINO]MLEICLKLVGCKSKKGLSSSSSCYLEEALQRPVASDFEPQGLSEAARWNSKE...[END_AMINO]

[START_I_SMILES]01[C@@H]([C@@H](O)[C@@H](O)[C@@H]1N2C(=O)NC(=O)C=C2)...[END_I_SMILES]

**Question:** What will be the docking score of this compound against the protein?

**Answer:** -8.8

**Figure 25: DockSTRING Format.** To construct the training set, we take the protein target and ligand sequences, pose a natural language question, and have the docking score as the answer.

# NLP Task Examples

- Encoder approach
  - Classify text snippets according to contents
  - Recognize relevant chemical separation entities of defined classes
  - Infer relationships between extracted entities

- Decoder approach
  - Directly query model with paper text as context, e.g.:
    - "*{context} Does this paragraph describe a chemical separation?*"
    - "*{cumulative-context} What extractant is used in this paper?*"
  - Relevance and quality of outputs will vary widely by default
    - Build and train re-ranking proxy reward functions



A 100 µg Th(IV) or 100 µg U(VI) was taken in 10 mL solution containing various concentrations of sulphuric acid ranging from 0.2 to 3.0 mol $L^{-1}$. The extractions were carried out in 125 mL of separating funnel. The aqueous acid solution containing metal ion was taken in 10 mL of 0.5–4.0% of 4-methyl-N–n-octylaniline in xylene for 4 min to separate the organic and aqueous phases. The metal loaded organic phase was stripped with strippant by wrist shaking of the funnel for 5 min. The 10 mL of 0.1 mol $L^{-1}$ nitric acid was used as strippant for Th(IV) while $2 \times 10$ mL acetate buffer of pH 4.5 was used as strippant in case of U(VI). The pH was adjusted with sodium hydroxide and acetic acid as per their requirement for the spectrophotometric determination. For quantification of Th(IV) xylenol orange was used as chromogenic agent while for quantification of U(VI) bromopyrogallol red was used [39].

P.S. More, *et al.* Extraction of Th(IV) and U(VI) with 4-methyl-N–n-octylaniline as an extracting agent. *J Radioanal Nucl Chem* **331**, 4149–4158 (2022).

Brookhaven
National Laboratory

# Data, Caveats, and Warnings

- Compiling first-round list of journals to construct training and validation data corpus
  - Emphasize open access journals (API-accessible, machine-parseable)
  - ML tools for parsing old PDFs and scanned papers, if necessary
- Bias and variety
  - e.g. underrepresented negative results
  - LLMs often pre-trained on general language, require domain adaptation
  - Data augmentation techniques may be necessary, depending on text volume
- Open questions and challenges
  - Maximizing value of expert annotation and validation effort
  - Verifying generated separations are practical

Brookhaven
National Laboratory

# Additional ML Opportunities

- Extracting knowledge from charts and tables
  - Useful for both recent and old papers



Fig. 1. Cation exchange distribution coefficients in HCl solutions. (Dowex 50-X4) (Ref.1)



Table II. Kd Values in H₂SO₄



Fig. 10 Extraction selectivity of metal ions from the simulated HLLW using the ionic liquids ([Bmim][NTf₂]: (a), [Bmim][NfO]: (b)) containing MPE-TDGA. ([MPE-TDGA]: 50 mM, Mixing time: 2 h, $V_{0,Aq}=V_{IL}$: 1 cm³)

F.Nelson,T.Murase and K.A.Kraus,J.Chromatog.,13,503(1964).
H. Oosugi, et al. "Extraction behaviors of platinum group metals from an aqueous HNO3 solution using ionic liquids containing a novel thiodiglycolamide-type extractant." *Journal of Radioanalytical and Nuclear Chemistry* 331.11 (2022): 4577-4585.

Brookhaven National Laboratory

# Additional ML Opportunities



| $I^{\pi}$ | $E_{\gamma}(E2)$ (keV) | $E_{\gamma}(E1)$ (keV) | $\lambda$ | $B(E2)/B(E1)$ ($\times 10^6$ fm$^2$) |
|---|---|---|---|---|
| | | Band 2 | | |
| $9^-$ | 376.3 | 970.5 | 0.020(4) | 3.16(63) |
| $11^-$ | 449.4 | 911.6 | 0.144(7) | 7.75(38) |
| $13^-$ | 517.9 | 868.5 | 0.43(3) | 9.86(69) |
| $15^-$ | 565.6 | 832.1 | 1.36(8) | 17.7(10) |
| $17^-$ | 611.3 | 807.9 | 2.8(2) | 22.6(16) |
| | | Band 4 | | |
| $6^-$ | 271.1 | 1127.8 | 0.121(8) | 154(10) |
| $8^-$ | 363.1 | 1045.7 | 0.639(8) | 151(2) |
| | | Band 4a | | |
| $12^-$ | 479.1 | 901.2 | 5.58(22) | 211(8) |
| $14^-$ | 491.5 | 790.7 | 20.5(10) | 460(22) |