# AI Under the Hood:
## What Natural Language Processing (NLP) Means for Pharma
### MODERNA, COHERUS MEDICAL AFFAIRS USE CASE

Matt Lewis[a], Richard Graves[b], Tim Mikhelashvili[c], Scott Navarro[a], Matt Booth[d], Ricky Li[a], Nga Tong[e], Paul Walden[f], Robert B Geller[f], Sandra Carrico[b], Walter Bender[b]

[a]MEDiSTRAVA, Boston, MA, USA; [b]Sorcero Inc., Washington, DC, USA; [c]Amedea Pharma, Inc., Westbury, NY, USA; [d]MEDiSTRAVA, London, UK; [e]Moderna Inc., Boston, MA, USA; [f]Coherus Biosciences, Redwood City, CA, USA

## OBJECTIVE

- We aimed to demonstrate the benefits of transforming literature monitoring from a manual report-based process into an intuitive, semi-automated, AI-driven process.

## BACKGROUND

- Intelligent Literature Monitoring (ILM) is a novel approach of augmenting and running continuous systematic literature searches around defined areas of therapeutic domains.

- ILM leverages Artificial Intelligence (AI) and its subset discipline Natural Language Processing (NLP).

- This enables the extraction of valuable insights from an expansive number of medical publications, as they are being published, into a centralized location.

  In 2018, the Google AI Language Division published BERT, "Bidirectional Encoder Representations from Transformers," a pre-trained algorithm capable of accurately gauging the context of each word in a sentence. BERT uses techniques including the ability to look forward and backward across passages to create an optimal representation of language. Once that representation is developed algorithmically, it can be leveraged for applications of Language Intelligence. BERT is Open Source and allows for permanent, continuous, crowd-sourcing, and machine learning improvement. The Sorcero LI platform leverages language models including BioBERT, while adding several proprietary algorithms to generate life sciences-specific solutions.

## METHODS

- We assessed two different approaches to continuous literature monitoring: (1) manual search and output, (2) semi-automated search, with AI integration and digital output, applied across search scenarios producing high (>250) and low (<100) volumes of results

- ILM was performed using the Sorcero LI Platform and the BioBERT language model against identical corpora

- We compared each approach for man hours, sensitivity, specificity, versatility/utility of outputs and depth of insights

- We also qualitatively assessed the digital output from a user experience perspective

## RESULTS

The ILM approach with integrated AI was the most favorable option for search strings with high volumes of search results (e.g. "oncology" or "COVID") (Figures 3 and 4).

- Compared with manually searching and reviewing the literature, this method resulted in a time reduction of 88-92%, along with 99.8% sensitivity and 95% specificity (Table 1).

- Sorcero's LI approach and continuous learning delivered double-digit absolute performance improvements across all study types, exceeding the 95% NPV threshold commonly accepted as that required for a regulatory-grade literature review solution (Table 2).

- For low-volume search results, (e.g. rare diseases) a digital tool improved accessibility, distribution, communication, archiving, and sorting/filtering.
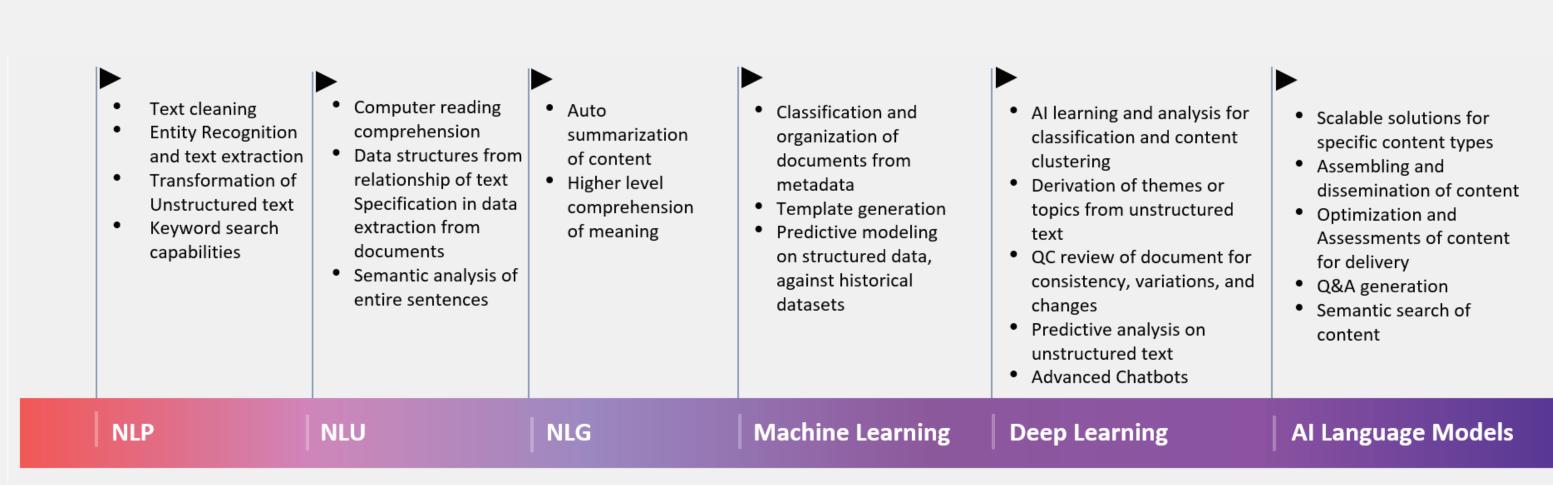


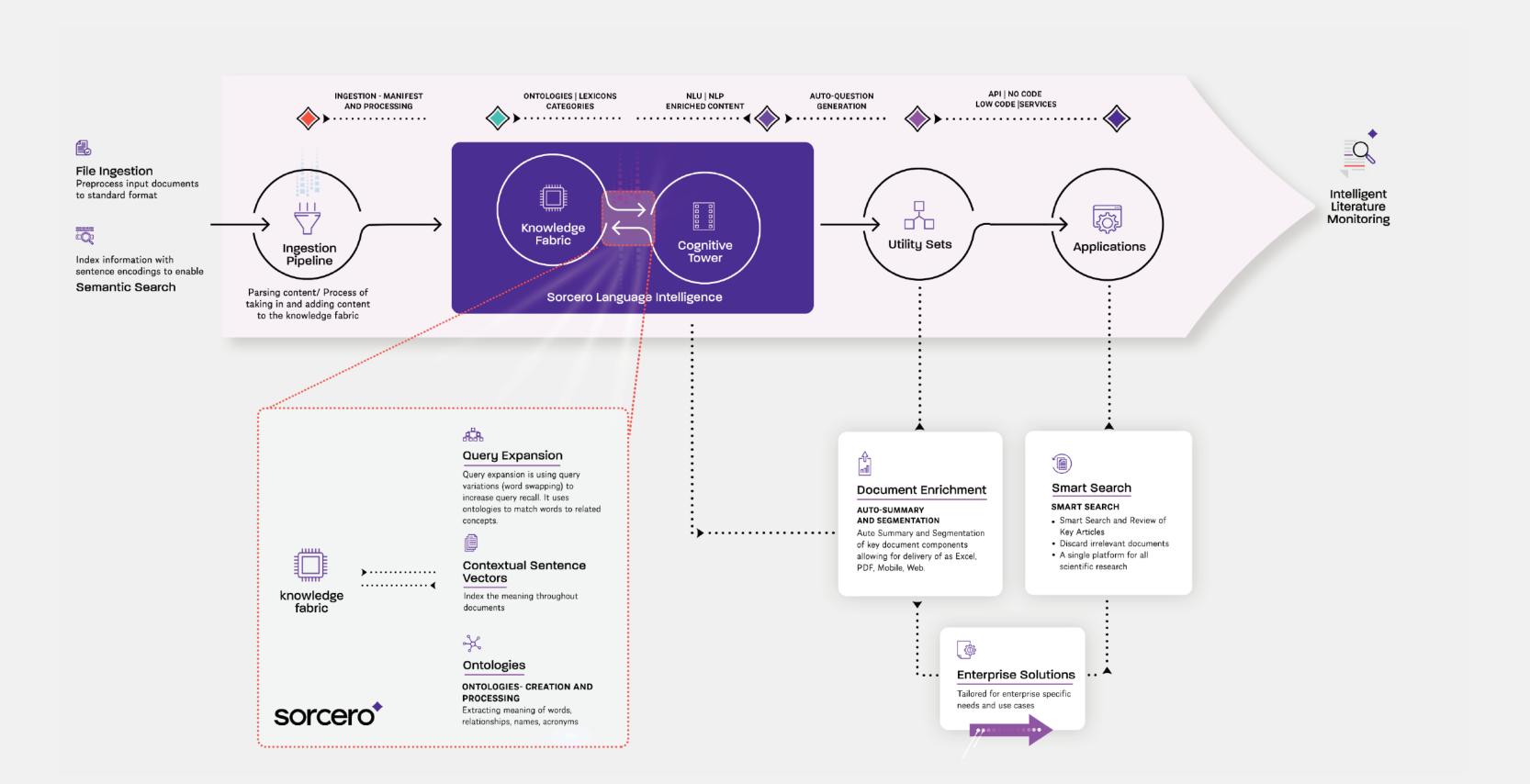Figure 1. Language Intelligent Capabilities spectrum
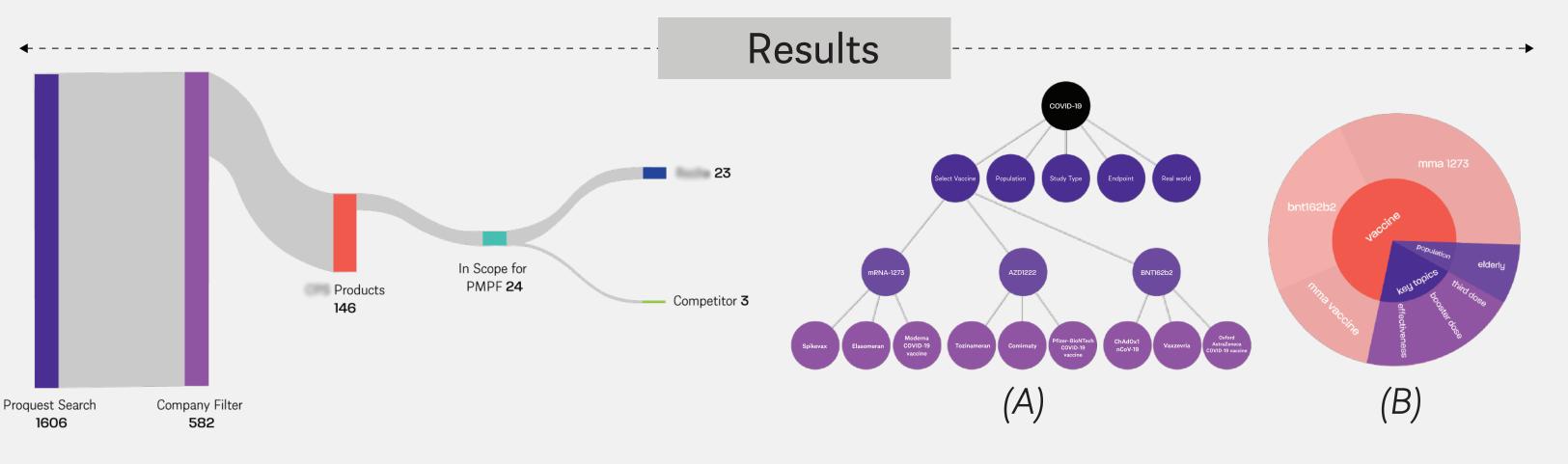


Figure 2. The ILM Approach with Integrated AI



Figure 3. Smart filtering through classification



Figure 4. Insight generation: auto-summary outputs from high volume ILM. A) Vaccine ontology tree; B) Keyword summary from mRNA vaccines dataset
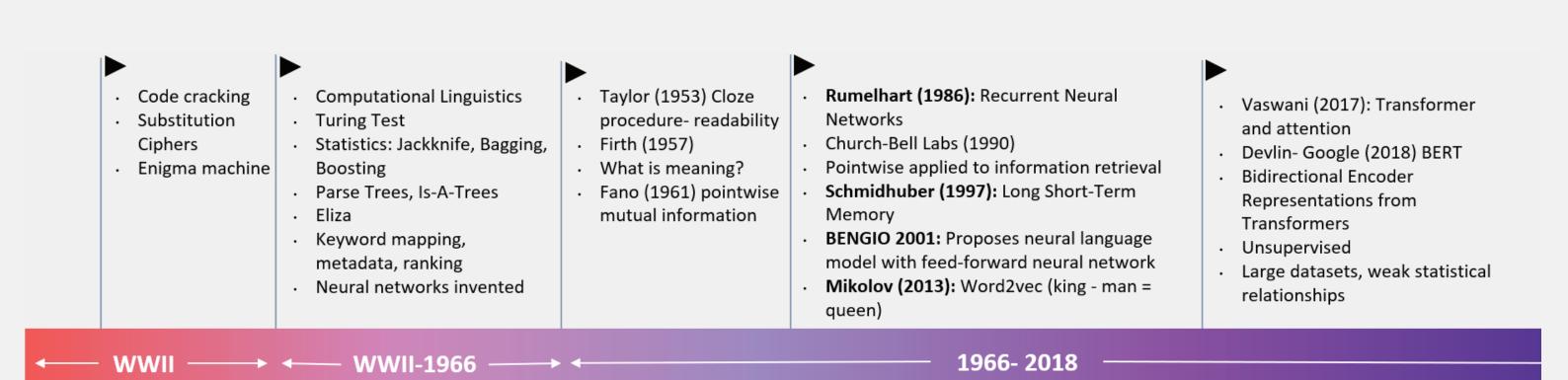


Figure 5. The Evolution of AI and Language Intelligence

| | Manual (Control Process) | ILM Process | Results |
|---|---|---|---|
| Total Articles | 1606 | 1606 | |
| # of Batches | 44 | 44 | |
| % Reviewed | 100% manual | 100% AI + Manual | |
| Company Criteria Relevancy Filter | 581 | 582 | 99.83% sensitivity |
| Product Relevancy Scored | 146 | 146 | |
| Relevancy Inclusion | 23 Human Determined | 23 AI Determined | >95% specificity |
| M.D./PhD. Time to Review | 715 hours (Avg. 27 mins/articles) | 55 hours Review + 25 hours of QA | 88% reduction in manual time |
| MD/PhD Cost to Review | $286,000 (@$400/hr) | $32,000 (@$400/hr) | $254,000 reduction in cost of time to review |

Table 1. Real-world results of the manual (control) process vs. ILM with integrated AI, performed for leading pharmaceutical enterprise medical affairs and publications teams

| | Analyte is subject | Performance study | Human study | Intended use |
|---|---|---|---|---|
| Bio BERT | 85.00% | 75.00% | 84.60% | 79.70% |
| Sorcero LI Platform | 96.50% | 88.50% | 96.90% | 91.10% |
| Differential | +11.50% | +13.50% | +12.30 | +11.40% |

Table 2. Comparative negative predictive value (NPV) from two AI literature monitoring platforms. Machine Learning and mathematical domain representation in the LI platform automatically tune general domain language models including Bio BERT for vastly higher performance in specific biomedical domains.

## CONCLUSIONS

**The ILM approach to systematic literature searching:**

1. Accelerates the current time-bound and episodic manual reviewprocess
2. Increases efficiency and productivity
3. Provides a deeper understanding of current complex, scientific content
4. Operates in near real-time
5. Allows automatic deployment of insights into a centralized and customized database

Digital platforms can compile and store results in a more versatile way, increasing functionality and creating efficiencies.

## DISCLOSURE

**Bibliography**

Lee, J. et al. Bioinformatics 2020;36:1234–1240
Vaswani, A et al. Proceedings of the Advances in Neural Information Processing Systems. 2017; 5998–6008

Sousa, D., Couto, F.M. Adv Inf Retr 2020;12036:367
Wei, Q. et al. Proceedings of the AMIA Annual Symposium. 2019;2019:1236
Zhu X. et al. Applied Sciences 2021;11:7318