# NLP Use Case Database Project

A bottom-up qualitative NLP success/failure database
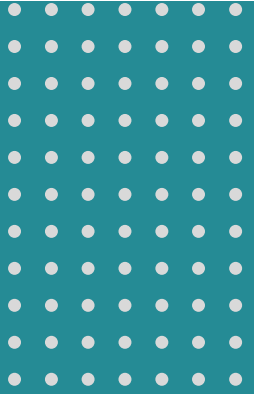
*Etzard Stolte, PhD, Roche*
*Bikalpa Neupane, PhD, Takeda*
*Birthe Nielsen, PhD, Pistoia Alliance*

*NLPdatabase@pistoiaalliance.org*

Pistoia Alliance

# Agenda

- The team

- The problem

- The project

- What has been achieved to date

- What is planned

- Get involved

Collaborate to Innovate Conference Boston 2022

# NLP Use Case Database

Project manager: Birthe Nielsen
nlpdatabase@pistoiaalliance.org

**Problem Statement:**

Pharma companies apply NLP methods in hopes of automation and insight generation. Although NLP algorithms have matured quite a bit during the past years, practical value for most NLP pilots tends to be poor, and very few NLP-driven projects are seen through to production. Exceptions are typically topics w/ good metadata quality, large training sets and willing business colleagues to verify results, and a serendipitous combination of technical expertise and suitable use cases.

**Value Proposition:**

This type of knowledge is of value to share in a pre-competitive manner among Pistoia Alliance members. A simple database could contain characterization of use case, data characterization, pipelines & algorithms used, quality criteria, outcomes, comments, etc. The value for participating members would be a reference database to help narrow down successful use case scenarios, less experimentation, and more successes.

**Project member roles:** Global Knowledge Management, Director of Advanced Technologies, Head of Text Analytics, Literature-based Scientific Discovery, Global Information and Analysis, NLP Lead/Engineer, Data Scientist

**Project champion:**
Etzard Stolte, Roche ; Bikalpa Neupane, Takeda.

**Steering committee (Currently forming):**

Roche     Takeda

$5-15K ensures a seat on the Steering Committee

**Project deliverables:**

- A bottom-up qualitative NLP Success Failure Database

- Agreed annotations of NLP use case methods and success/failure criteria

- Collaborative insight into why NLP use cases may fail or succeed with an industry –wide view

| Q1 2022 | Q2 2022 | Q3 2022 | Q4 2022 | 2023 |
|---|---|---|---|---|
| Roundtable Discussion | Presentation of Use Cases | **Draft Database Created, First Use Cases Added** | Populate 50-100 Use Cases | Standardize nomenclature, continue populating, analyze success prediction criteria |

https://drive.google.com/file/d/1AkmetDVgCH92zBJ4HZjuVhSgn5H1XqZU/view?usp=sharing

# Project timeline till date

Roche ideas proposal
Round table discussion with interested pharma members

Draft Use Case Database

First Use cases added

Working group meeting to agree on annotation of NLP use case methods and success/failure criteria

Collection of use cases to be completed (50-100 Use cases)

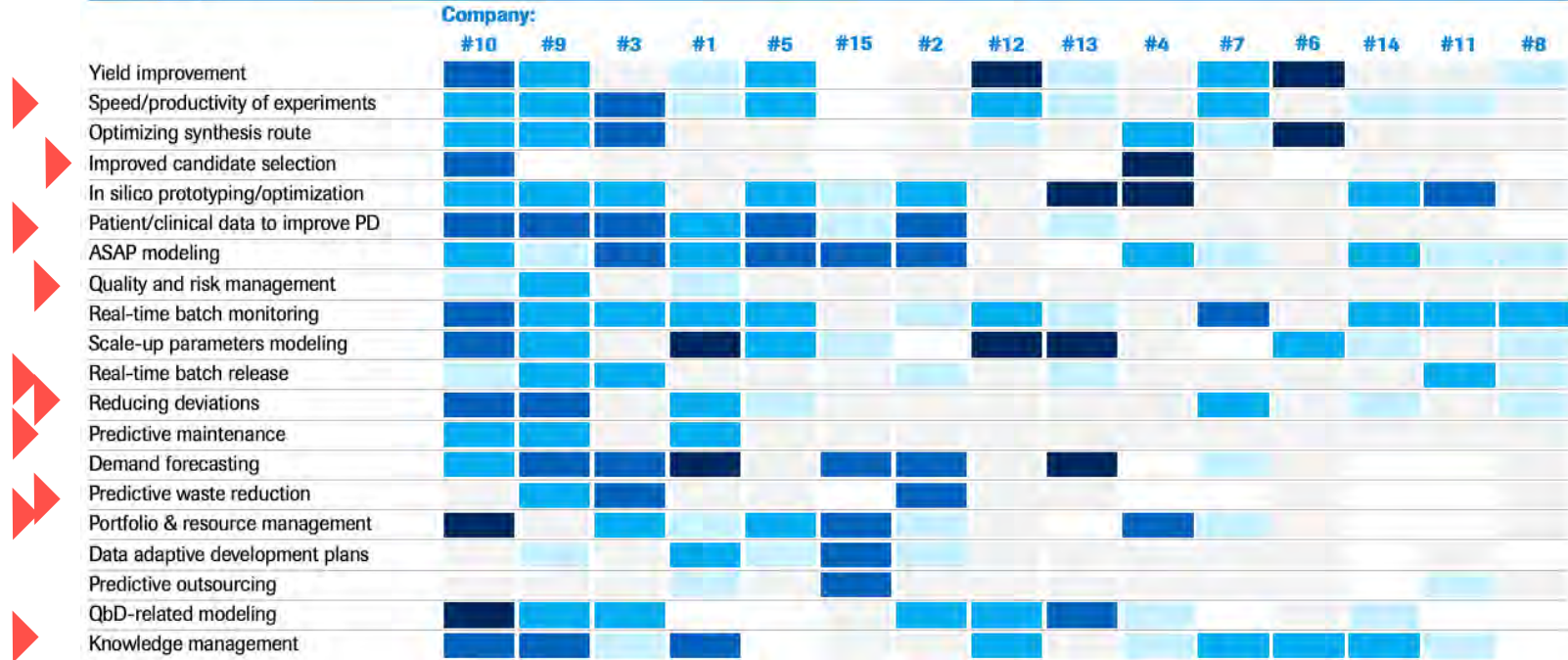Collaborative insight to why NLP use cases may fail of succeed – wider view.
Next steps discussion

| Q1 2022 | Q3 2022 | Q3 2022 | Q4 2022 | Q1 2023 |

# Machine Learning (NLP, NER) uses cases in Manufacturing **theory**

## Wide range of use cases



What are the advanced analytics use cases employed by your CMC organization? Mature | Early | Advanced | Planned | Not planned | No visibility
Status of adoption

# Hype Cycle for AI, ML, NLP

*Gartner, 2020: real world applications of AI/ML/NLP struggle with quality*

reality



© Pistoia Alliance

# Natural Language Processing (NLP)

*reality*

*Challenges ? Hell yes - it simply does not deliver the quality we need in Pharma*

- Pharma is moving beyond buzzwords to value generation
- NLP is more of an art, then engineering

- Observations
  - Hard to determine up front, if a use-case will work
  - "successful" pipelines do not work for slightly different data
  - Manual curation, training sets are exceedingly work-intensive
  - Success seems very use case specific
  - Expensive to find talent that actually delivers the quality we need

# Pistoia Alliance - NLP Expert Group

*Value Proposition*

- Share best-practice – get an idea what works & what does not work
- Build database of use cases to be shared in this group; incl. failures

- Obvious challenges to address

  - Can we come up with a categorization that makes sense?

  - If NLP is an art, can we structure use-cases in a meaningful way?

  - Is a phenotypic classification sufficient (aka like Linne)?

  - Can we deliver on the effort to identify and enter your use cases iteratively into some database?

*Systema naturæ sistens regna tria naturæ, in classes et ordines, genera et species redacta tabulisque æneis illustrate, Systema Naturae, 1768, Carl von Linné*

# Natural Language Processing Landscape

© Pistoia Alliance

# Meaning Based Computing *use cases*

## Use Cases types and examples

| | Concepts | Trend Detection | Quality Controls | Risk Prediction | Cause Analysis |
|---|---|---|---|---|---|
| **Type** | *Automatically extract and correlate concepts using semantic tools* | *Interactively explore emerging concepts / trends* | *Monitor operational & content quality measures hidden in un-structured information* | *Real-time classifiers to sense risks across sources, sites & organizations* | *Extract relationships among concepts in a Knowledge Graph* |
| **Examples** | • find documents/experts<br>• Automatically add quality tags to lab analysis<br>• Show potentially recurring deviations for new TrackWise entries | • sense patterns in health authority records<br>• early warning signs for batch or release records<br>• off-label use of products | • detect mis-filed documents in submission structures<br>• flag spurious relationships in manufacturing logs | • sentiment analysis of HA ePARs with respect to submission success<br>• emerging issues in batch records | • understand cause / effect relationships in manufacturing deviations<br>• visualise relationships / dependencies inside SOPs |

# The Use Case Database

Contributions till date:

# Where our current database stands today

- Named-entity recognition (NER)
- Concept Extraction (Taxonomy extraction, concept tagging, document classification)
- Text-Classification
- Unsupervised topic modeling
- Extractive auto-summarization
- Similarity Search
- Chatbots/Virtual Assistants

# Database content

- Use case
  - Deployment type
  - Project level
  - Company
  - Language
  - etc
- Data
  - Description
  - Size
  - Issues
  - Preparation
  - FTE involved/team
  - etc
- Algorithm
  - Summary
  - Version
  - Library
  - etc

- Outcome
  - Learnings
  - Success achieved
  - etc
- Discussions
  - Effort/ benefit ratio
  - Non-NLP approach
  - Success
  - Support
- Outlook
  - Build vs buy recommendation
  - Other applications of technology
  - etc

# Pistoia Alliance - NLP Expert Group   worry

## What do we **NOT** want to do

- Encyclopedia of NLP methods

- Another NLP discussion group

- Vendor messages, vision documents, general statements

- Evangelism of a particular NLP school of thought or technology

- ???



https://mattturck.com/, November 2021 – Machine Learning & Data Landscape

# NLP Build vs Buy

- ❑ *Language is hard – ambiguous,* Challenges because human expression is diverse, often ambiguous, affected by age, demographics, socio-economics, medium/channel & regional attributes of the author.

- ❑ Building a POC NLP model is trivial. However, gathering domain specific dataset, training large models and deploying them reliability (serving thousands of requests at low latency) is extremely hard.

- ❑ Running Experiments in Parallel and Serving Models is computationally expensive

- ❑ Market Barriers – Lack of Clean and accurate data -   Lack of human-experts;

- ❑ Compliance is not trivial (PII, PHI, GDPR, FERPAA) – Customizability, AI Explainability, and Fairness are daunting challenges to fulfill without a large, in-house data science team

- ❑ Head Count + Maintenance Costs  + Infrastructure Costs ➜ ROI ?

# What if you were to purchase, checklist

- ❑ Return on Investment (ROI)
- ❑ Domain Fit Accuracy of the results
- ❑ Production Readiness - Scalability and Performance, Deployment Options
- ❑ Interoperability and Fit with existing software stack
- ❑ Integration with existing ML pipeline
- ❑ Flexibility - NLP custom models
- ❑ Usability
- ❑ Industry Leaders – Look at who is doing the real NLP work
- ❑ Compare and contrast regular use cases

# We are eager to have you….

**More information here:**

- https://www.pistoiaalliance.org/projects/current-projects/natural-language-processing-use-case-database/

**Get in touch:**

- nlpdatabase@pistoiaalliance.org

## Questions

nlpdatabase@pistoiaalliance.org