

# The Automatic Scientist will be a Data System

Stratos Idreos

Harvard University

## 1. VISION: AUTOMATIC SCIENTIST

For thousands of years science happens in a rather manual way. Mathematics, engineering and computer science provide the means to automate some of the laborious tasks that have to do with computation, data collection and management, and to some degree predictability. As scientific fields grow more mature, though, and scientists over-specialize a new problem appears that has to do with the core of the scientific process rather with the supporting steps.

*It becomes increasingly harder to be aware of all research concepts and techniques that may apply to a given problem.*

A critical step when generating new concepts and explaining existing phenomena is being aware of existing concepts/techniques and how they can be combined. This is true in a single scientific area at a time but it is exacerbated across areas, i.e., we can easily miss opportunities to leverage advances across unfamiliar areas. In addition, this is important when we try to explain a phenomena, e.g., in medicine when a doctor tries to match an existing patient to the vast space of possible conditions or when a surgeon tries to decide what the best next step should be in real time during a surgery. Computer science already helps in significant ways. Recent advances include automation in feature engineering [1] and hypothesis generation [2] by analyzing past research. What we propose here is to take one step further and map existing research concepts (at a fine granularity) in a unified model and then treat that as a data management and analytics problem.

We envision a future when scientists rely on software tools not only to collect, manage and analyze data but also to generate ideas, make it easy to solve problems by combining existing research concepts across one or more areas, get suggestions about possible next steps, and find errors, semi-automating the scientific process. In the same way that modern data systems organize the world's data and facilitate data processing in numerous areas in businesses and sciences, we envision that "auto-science data systems" will be able to abstract and speed up the scientific process across numerous scientific fields. Auto-science systems are "scientist-in-the-loop" environments that leverage intuition and experience of scientists along with automation provided by the tools for the parts of the scientific process that rely on combining/tuning past work in novel ways.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

## 2. RESEARCH PATH

We argue that auto-scientist tools will be data systems, largely relying on similar procedures and properties that have been pioneered in the data systems community.

**Models.** Similar to the relational model, we need a unified model that is generic enough to capture the necessary information needed to describe basic research concepts and their combinations across sciences. Such a model should capture research concept characteristics, such as when it applies, its dependencies, cost metrics and side-effects, and tuning parameters.

**Declarative Research Processing.** Once we have a model, we can develop declarative languages over this model to provide browsing and analytics over research concepts. These new class of queries will be tailored to speed up research by leveraging the structure of research concepts e.g., find how we can best solve problem X (specific requirements  $R$ ) given existing research concepts that can be combined to satisfy  $R$  or find the weakest link in a research solution that can be described as a combination of concepts.

**Modularity & Optimization.** Modularity is critical for scalability and to be able to utilize fine-grained combinations of research concepts. Using experience from modular systems with clean APIs, auto-scientist systems will be easy to extend with new concepts and to synthesize research solutions. Then any algorithm, data structure developed to enhance browsing of scientific concepts over the generic infrastructure will apply universally across all scientific fields and similarly to how modern systems optimize queries right algorithms/data structures should be abstract to the scientists.

**Summary.** The auto-scientist is a data system that stores research concepts and allows scientists to interactively navigate the research space to speed up research. It also brings opportunities in education of young scientists and verifiability of research solutions. It resembles in many ways the path that the database community followed to develop systems for data management - only this time this is an effort that has to happen in collaboration with the target scientific fields to develop the right models and expressibility. In addition, we expect many research challenges to be analogous to database systems, e.g.: Can we indeed have a single model or we need multiple models which brings data integration problems? Do we need polystores that store concepts from different fields? Our community has several existing ideas to seed such a research landscape.

## 3. REFERENCES

- [1] M. R. Anderson, D. Antenucci, V. Bittorf, M. Burgess, M. J. Cafarella, A. Kumar, F. Niu, Y. Park, C. Ré, and C. Zhang. Brainwash: A Data System for Feature Engineering. In *CIDR*, 2013.
- [2] M. Nagarajan, A. D. Wilkins, B. J. Bachman, I. B. Novikov, S. Bao, P. J. Haas, M. E. Terrón-Díaz, S. Bhatia, A. K. Adikesavan, J. J. Labrie, S. Regenbogen, C. M. Buchovecky, C. R. Pickering, L. Kato, A. M. Lisewski, A. Lelescu, H. Zhang, S. Boyer, G. Weber, Y. Chen, L. A. Donehower, W. S. Spangler, and O. Lichtarge. Predicting Future Scientific Discoveries Based on a Networked Analysis of the Past Literature. In *SIGKDD*, 2015.