# ROBOKOP and the Biomedical Data Translator: Efficiently Leveraging a Distributed Data Ecosystem

*Jim Balhoff[1], Chris Bizon[1], Steve Cox[1], **Karamarie Fecho**[1,2], Yaphet Kebede[1], Kenneth Morton[3], Alexander Tropsha[1,4], **Patrick Wang**.[3]*
*[1]Renaissance Computing Institute (RENCI), University of North Carolina at Chapel Hill*
*[2]Copperline Professional Solutions*
*[3]CoVar Applied Technologies*
*[4]School of Pharmacy, University of North Carolina at Chapel Hill*

# NIH Mission...*is to seek fundamental knowledge about the nature and behavior of living systems and the application of that knowledge to enhance health, lengthen life, and reduce illness and disability*
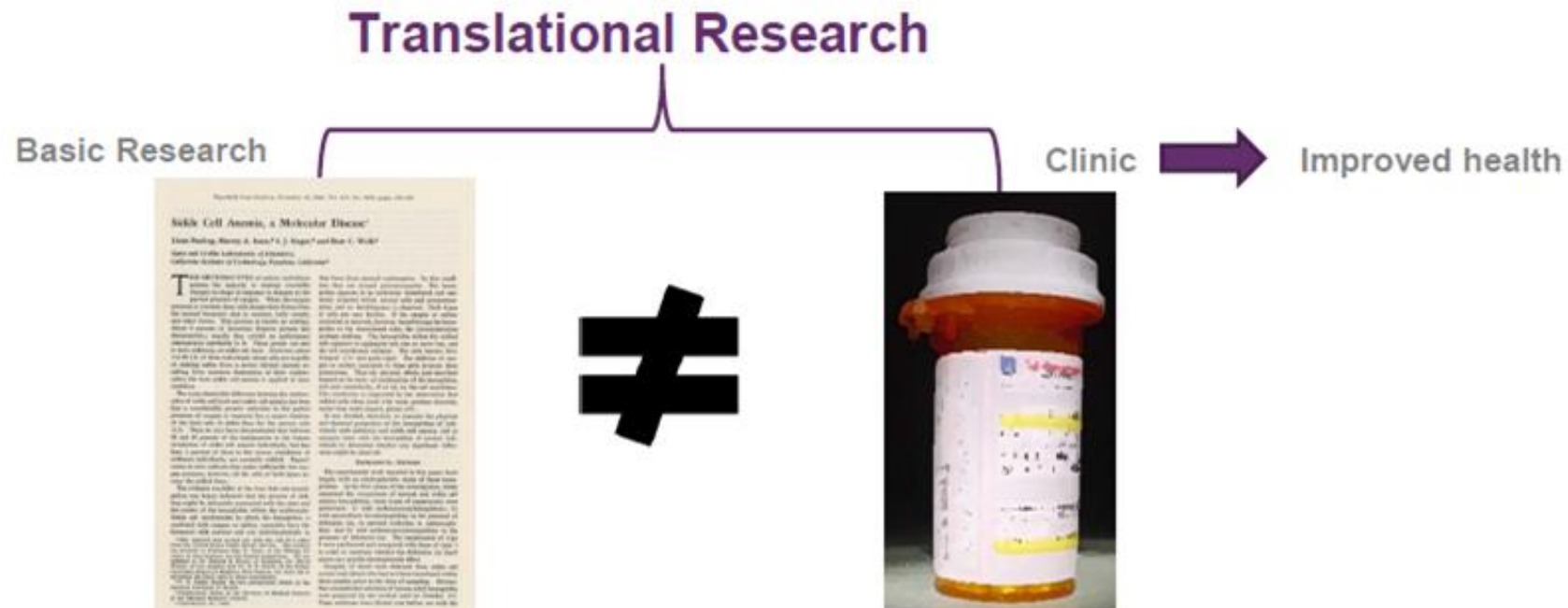
# NCATS Mission...*is to catalyze the generation of innovative methods and technologies that will enhance the development, testing and implementation of diagnostics and therapeutics across a wide range of human diseases and conditions*
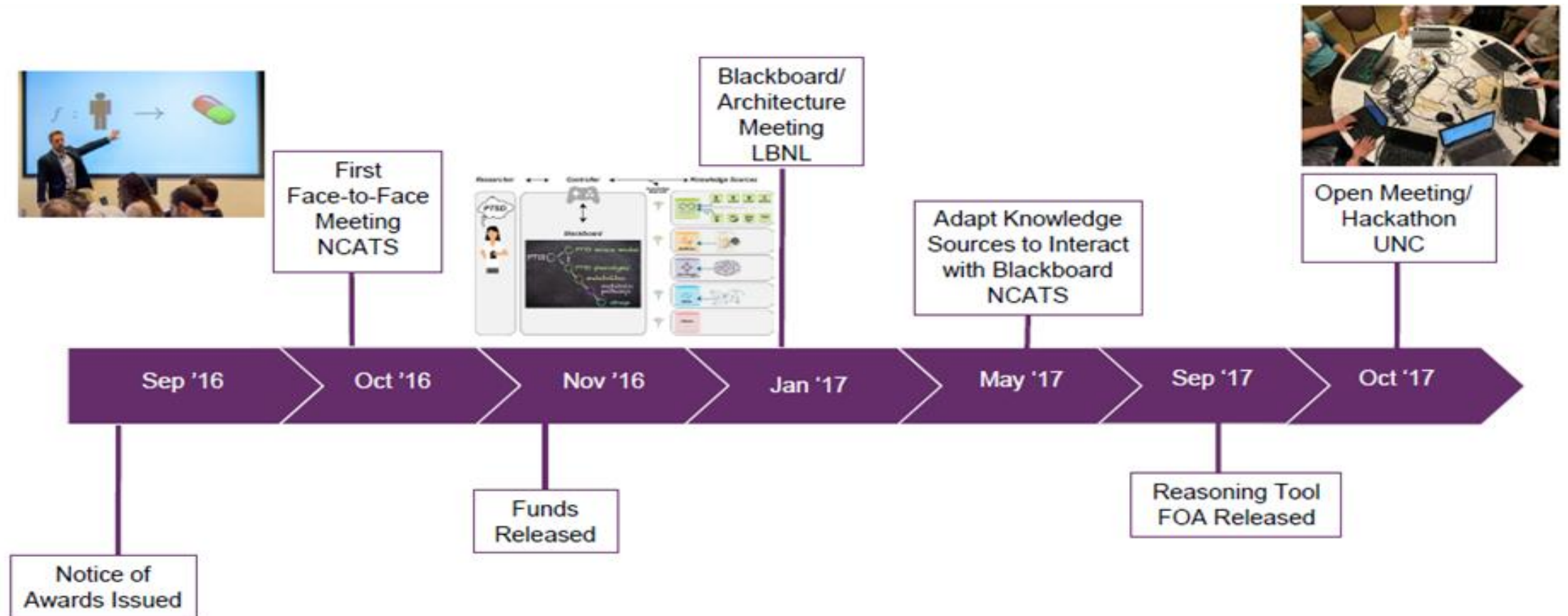


*Image created by NCATS*

# NCATS Biomedical Data Translator Program

- Overarching goal is to promote "**serendipitous**" discovery
- Primary goal is to build infrastructure to support and facilitate **data-driven** translational research on a large scale

- Essential aim is to **link as many datasets as possible** with one another and allow them to be **cross-queried** and **reasoned over** by translational researchers…

    …at **speed**, with **minimal barriers**, **scalability**, **accuracy**

- Fundamental tenet of the program is **open data** and **open collaborative team science**
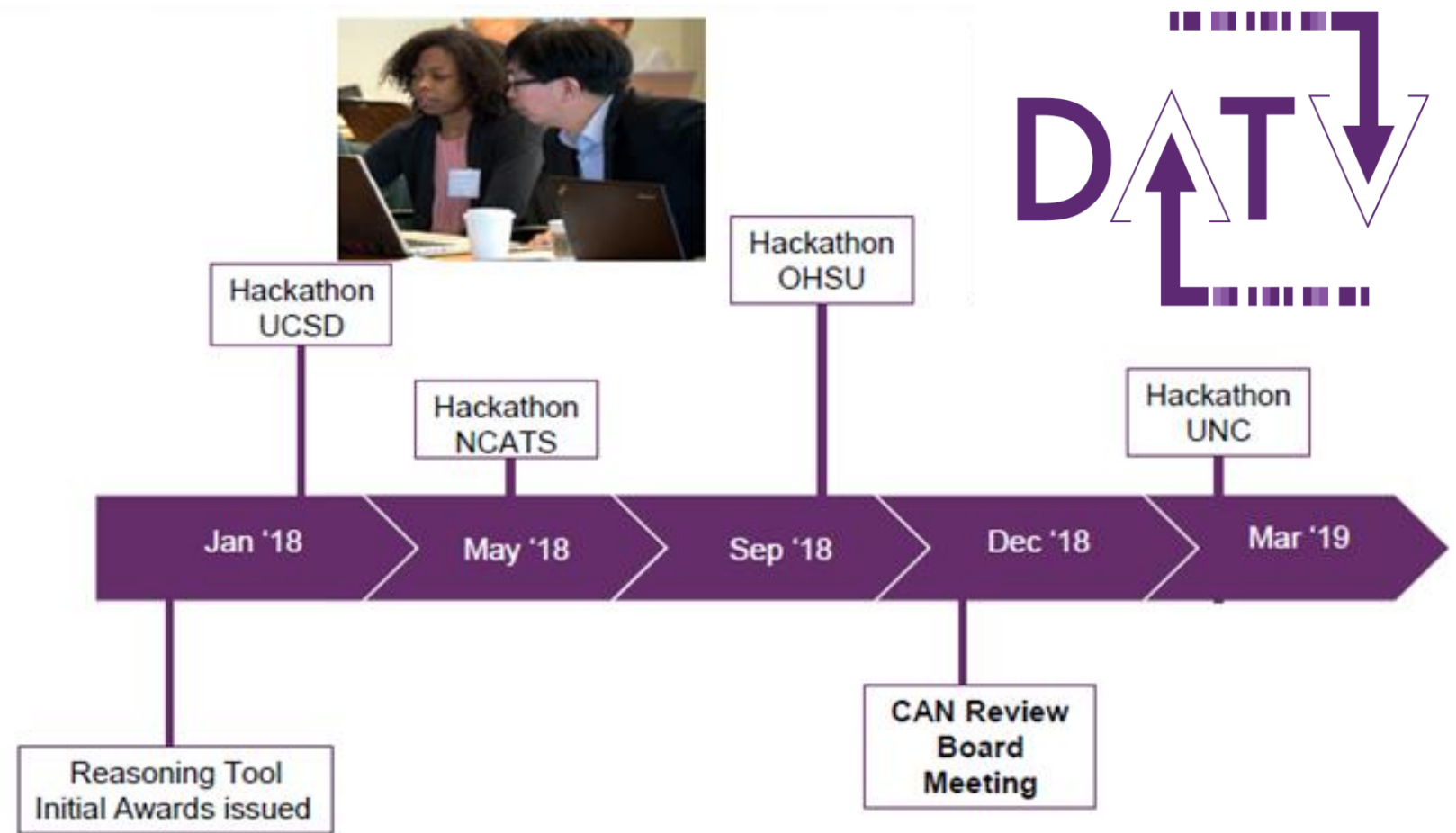
# NCATS Biomedical Data Translator Program
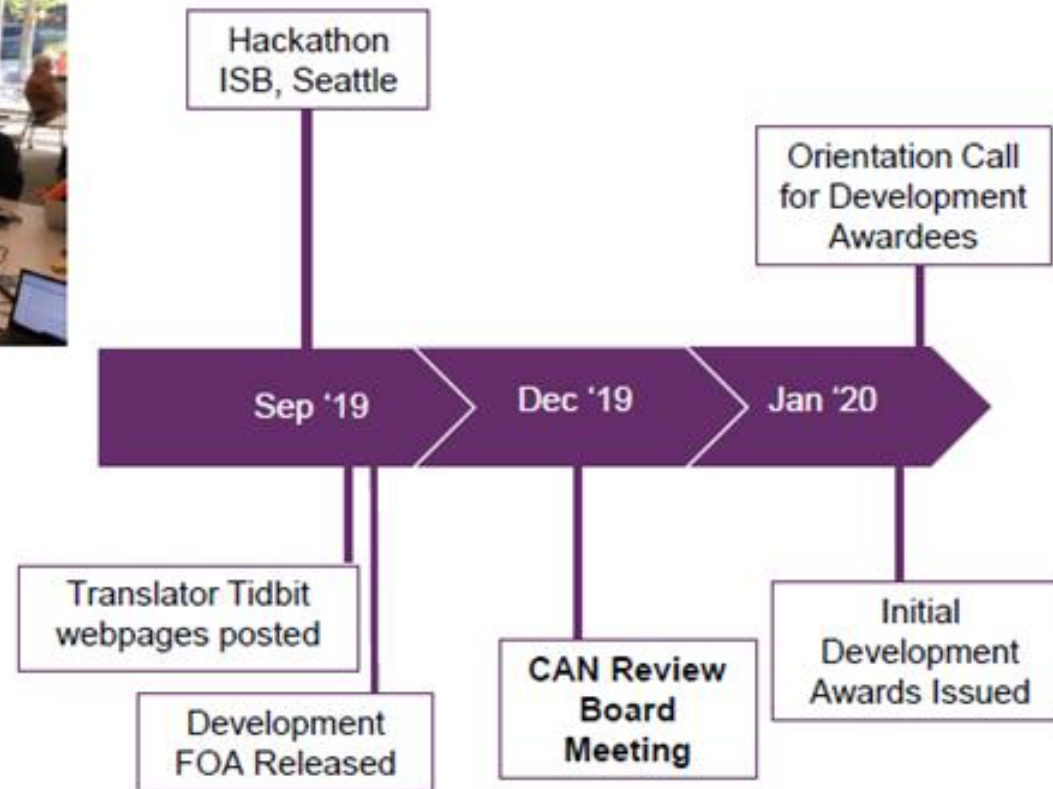## *'Other Transaction Award (OTA)'*



Blackboard/ Architecture Meeting LBNL

First Face-to-Face Meeting NCATS

Adapt Knowledge Sources to Interact with Blackboard NCATS

Open Meeting/ Hackathon UNC

| Sep '16 | Oct '16 | Nov '16 | Jan '17 | May '17 | Sep '17 | Oct '17 |

Notice of Awards Issued

Funds Released

Reasoning Tool FOA Released

*Image created by NCATS*

# NCATS Biomedical Data Translator Program



Hackathon UCSD

Hackathon NCATS

Hackathon OHSU

Hackathon UNC

Jan '18    May '18    Sep '18    Dec '18    Mar '19

Reasoning Tool Initial Awards issued

CAN Review Board Meeting

*Image created by NCATS and adapted here*

# NCATS Biomedical Data Translator Program



Hackathon
ISB, Seattle

Orientation Call
for Development
Awardees

Sep '19    Dec '19    Jan '20

Translator Tidbit
webpages posted

Development
FOA Released

CAN Review
Board
Meeting

Initial
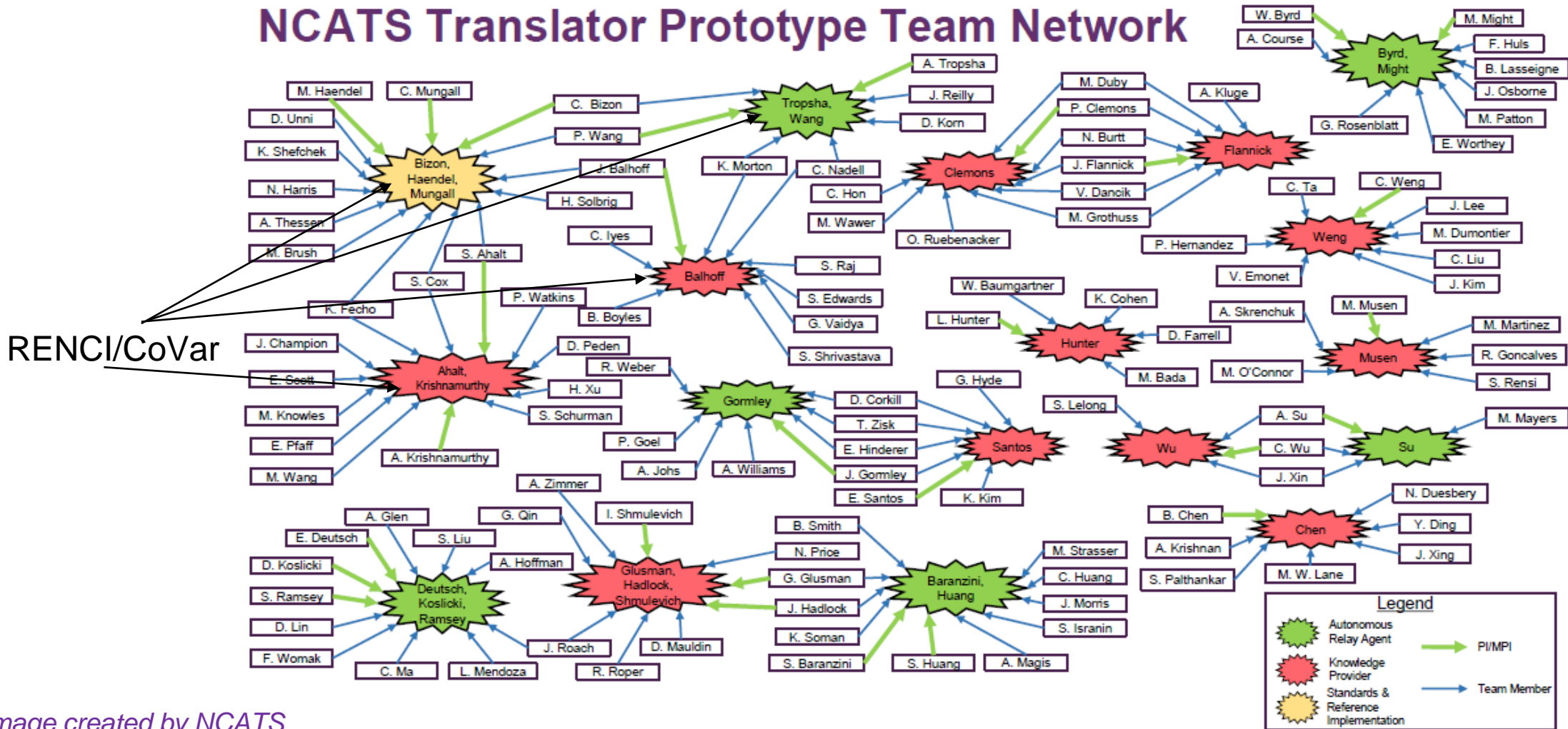Development
Awards Issued

*Image created by NCATS*

# The Biomedical Data Translator Consortium:
# 11 Teams, 28 institutions, ~200 team members
## Phase I: 09/24/2016 - 12/31/2019

# The Biomedical Data Translator Consortium:
## Phase II: kick-off 01/27/2020



Image created by NCATS

# Translator Vision



"Two hundred years ago, chemists created a comprehensive enumeration of the elements and systematic relationships among them. We envision the Translator doing the same for translational science."
— Christopher P. Austin, MD, director of NCATS, with Christine M. Colvis, PhD, Noel T. Southall, PhD

*Image created by Julie McMurry, with input from the Biomedical Data Translator Consortium*
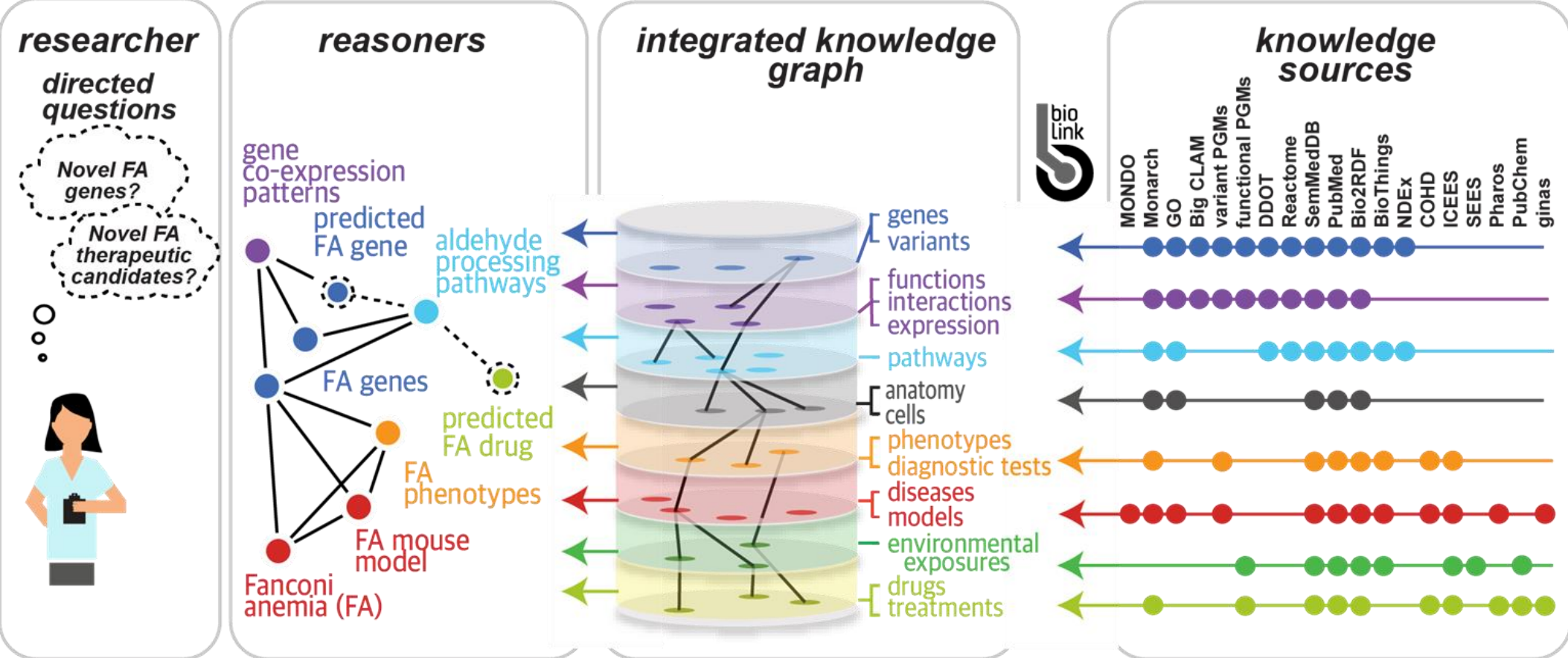
# Prototype Translator System Architecture



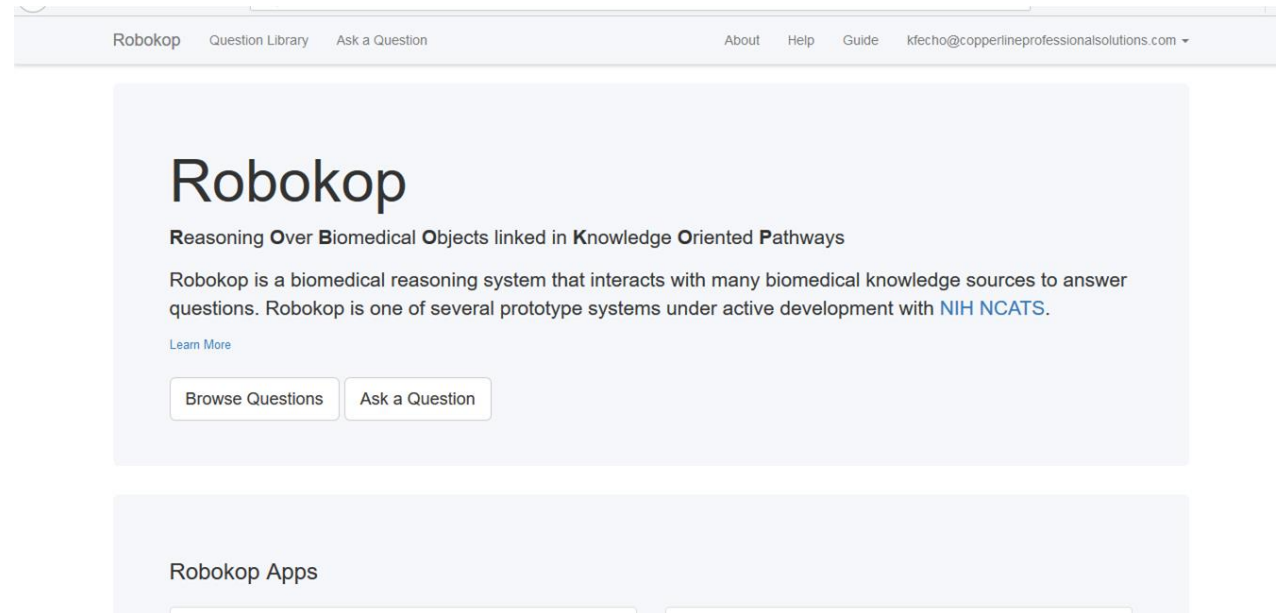*Image created by Julie McMurry, with input from the Biomedical Data Translator Consortium*

# ROBOKOP (Reasoning Over Biomedical Objects linked in Knowledge-Oriented Pathways)

- An open knowledge graph (KG)–based question-answering system
- Can be used to derive mechanistic insights into real-world observations

*what biological pathways might explain the association between exposure to carbon monoxide and multiple sclerosis?*

*how does hair dye influence susceptibility to breast cancer?*

*what genes might explain the effectiveness of isopropyl alcohol in the treatment of cyclic vomiting syndrome?*



robokop.renci.org

# Demo

## Ask a Question

**Question Title**

What genetic conditions may provide protection from Ebola?

Machine Question Editor - Question Graph ⓧ

⊞ Add Node
✚ Add Edge

0: Ebola hemorrhagic fever ——— 1: Gene ——— 2: Genetic Condition

Advanced

Reset 🗑    Download 📥    Submit ➤

# Services

- "manager"
  - manages users
  - stores queries and results
  - presents custom web GUI
    - building and running queries
    - exploring results
- "ranker" / "messenger"
  - receives queries from public API, in JSON format
  - transpiles into Cypher query
  - collates, scores, and returns results

- "builder" / "interfaces"
  - consults registry of (third-party) APIs
  - identifies sequences of queries that could combine to produce answers
  - iteratively queries to gather relevant data
  - stores data in Neo4j database "robokopdb"
- "robokopDB"
  - publicly available Neo4j browser

# Builder/Interfaces

# Ranker

# ROBOKOP KG

- > 6M nodes and 140M edges
- ~30 underlying biomedical data sources and bio-ontologies (e.g., KEGG, Monarch services, DrugCentral, Pharos, Comparative Toxicogenomics Database, Monarch Disease Ontology, Gene Ontology, ChEBI, Human Phenotype Ontology)
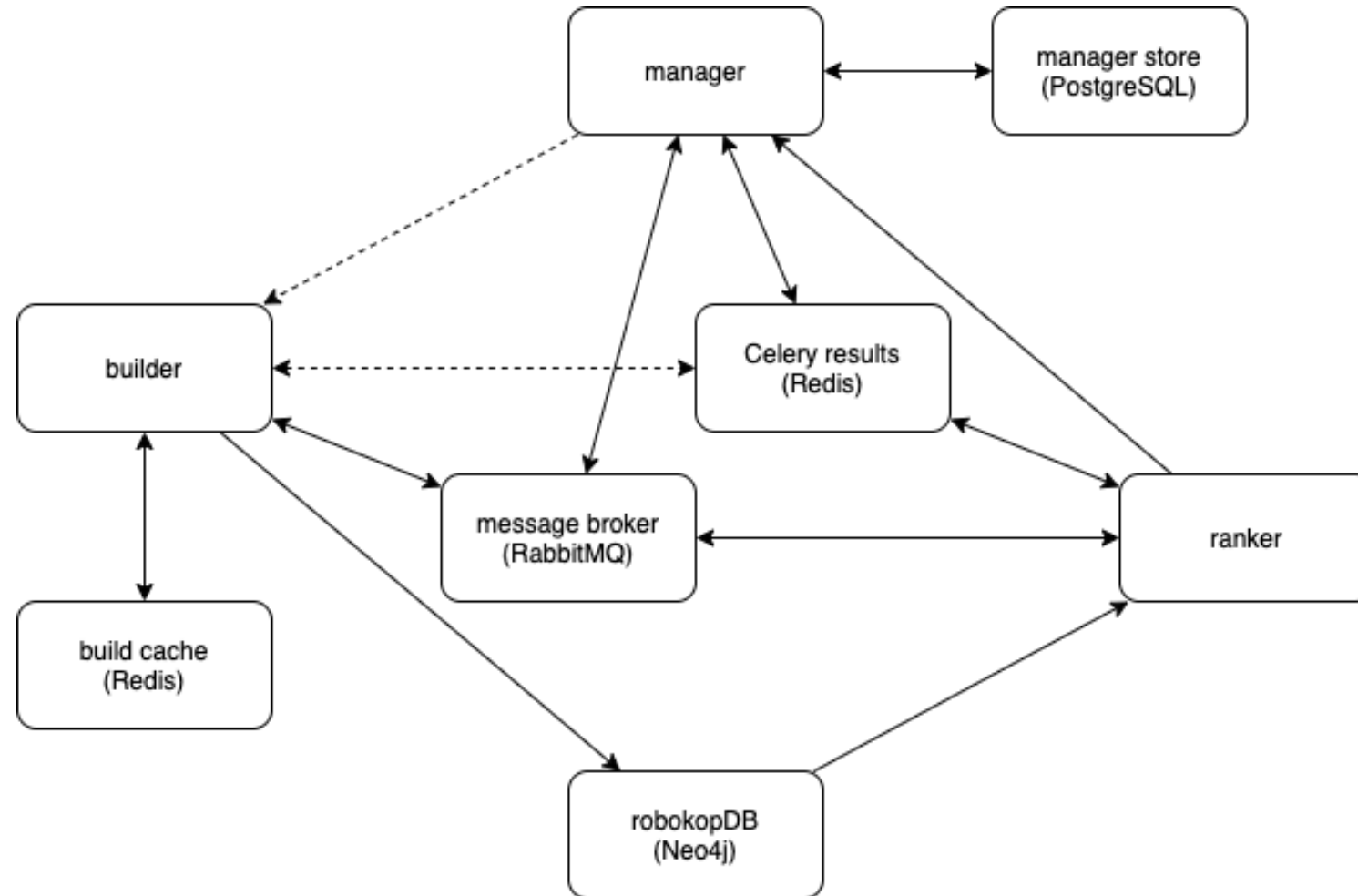- Extensible to non-biomedical domains and knowledge sources



robokopkg.renci.org

# ROBOKOP Tooling

- Supervisor: process control system
  - e.g. Flask process, Celery process, …
- Celery: distributed task queue
  - master Celery (Python) process
  - multiple Python worker processes
  - jobs distributed to workers via message broker (RabbitMQ)
  - results stored by workers in results database (Redis)

- NGINX, Gunicorn, Flask: web server
  - NGINX: reverse proxy
  - Gunicorn: WSGI HTTP server
    - multiple worker processes
  - Flask: web server framework
- Redis, Neo4j, PostgreSQL: caching
  - Third-party API calls (Redis)
  - Biomedical data (Neo4j)
  - Entire questions/answers (PostgreSQL)
- SQLAlchemy, Postgraphile + GraphQL
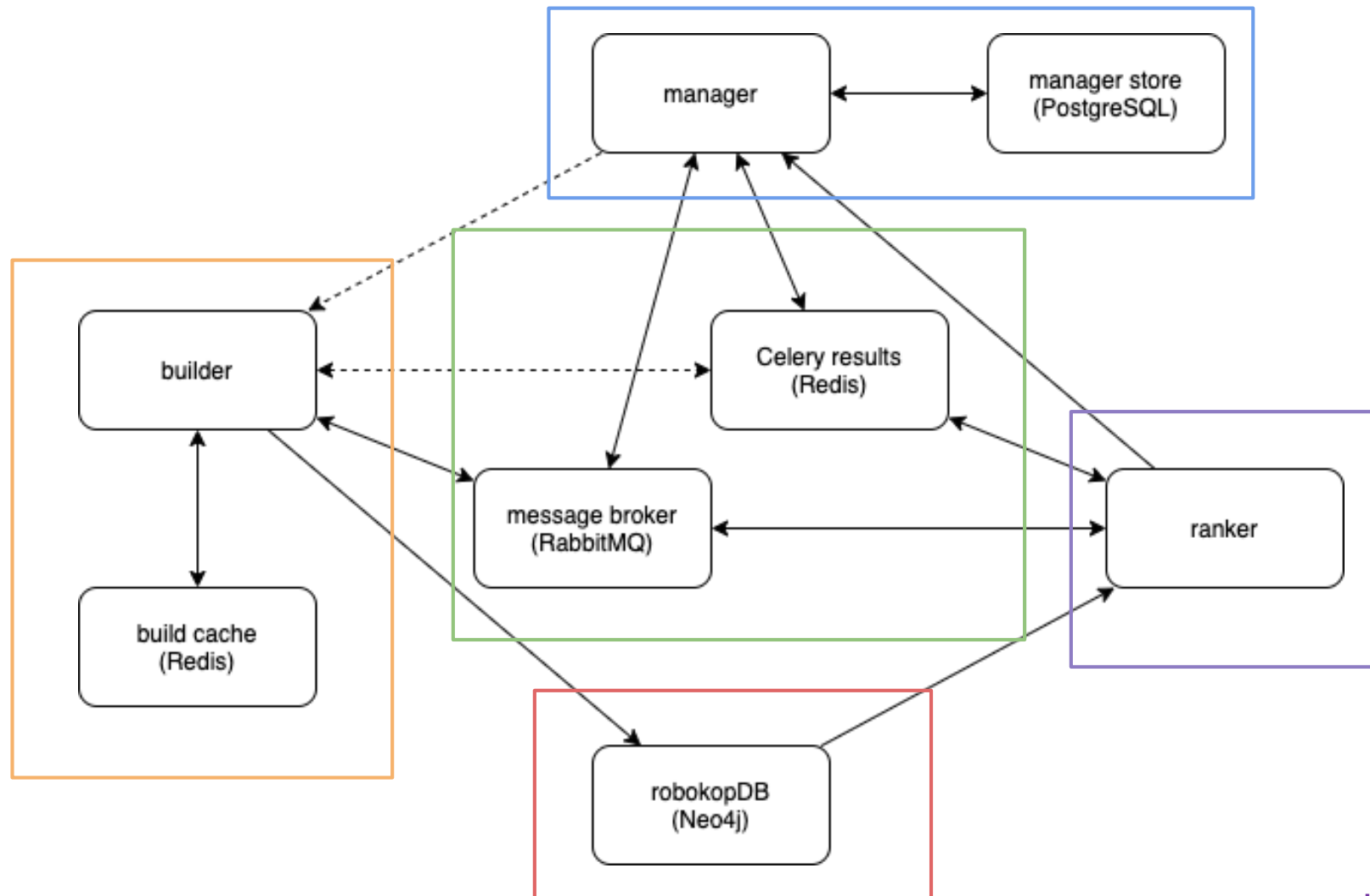  - interface with PostgreSQL

# Containers

- manager
  - Python/Flask
- ranker/messenger
  - Python/Flask
- builder/interfaces
  - Python/Flask
- robokopDB
  - Neo4j

- reverse proxy
  - NGINX
- build cache
  - Redis
- manager store
  - PostgreSQL
- Celery results
  - Redis
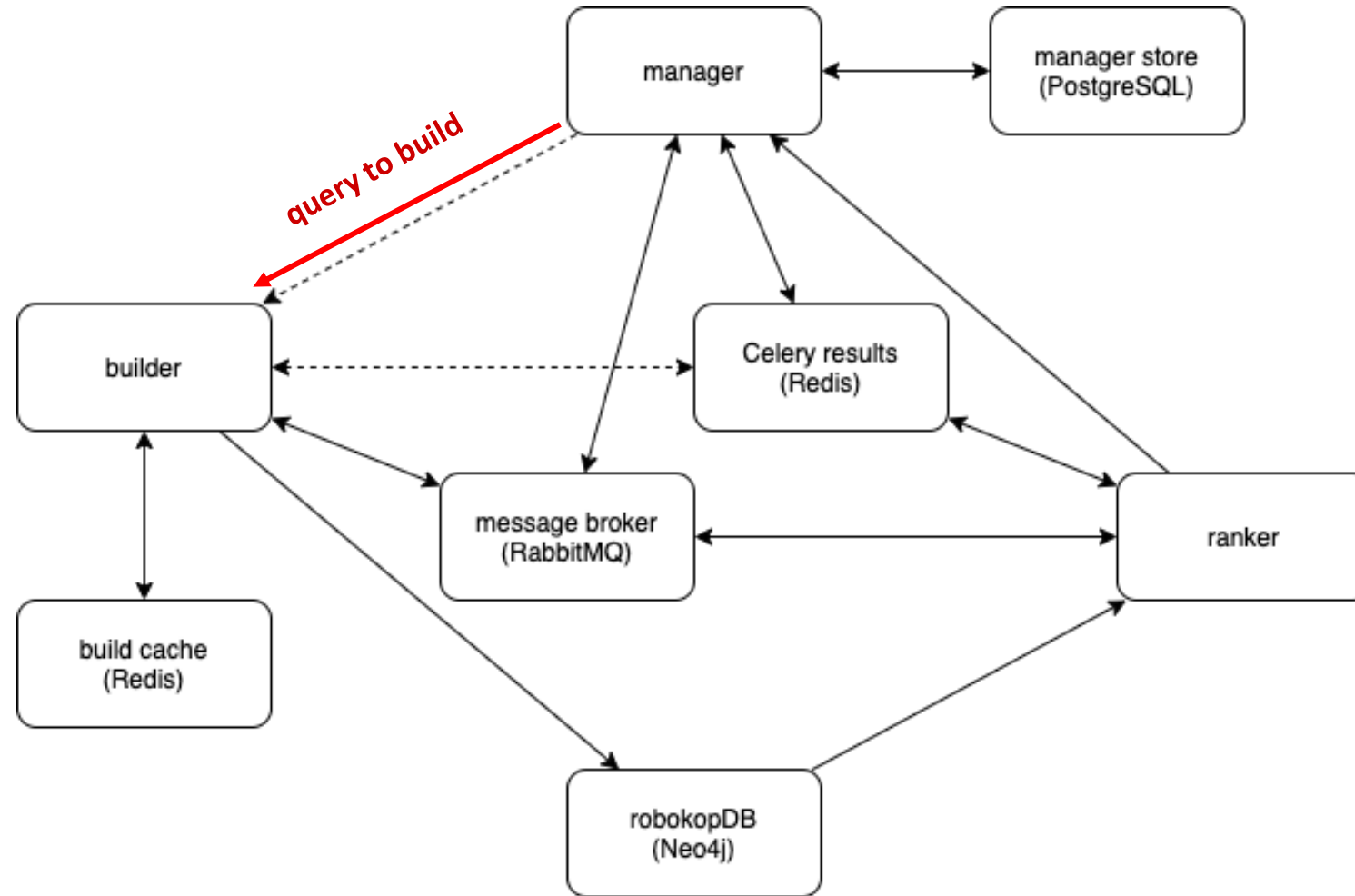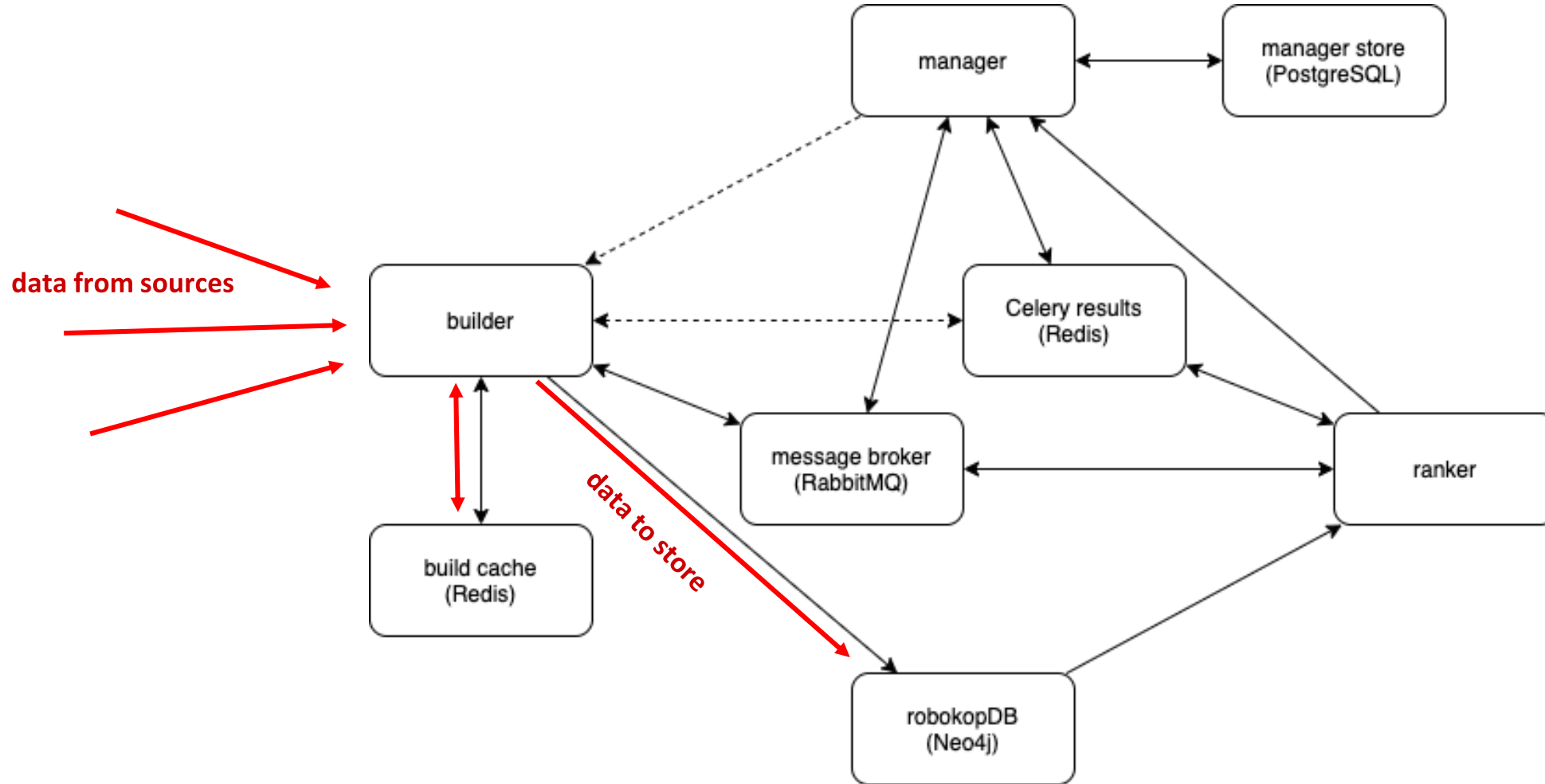- message broker
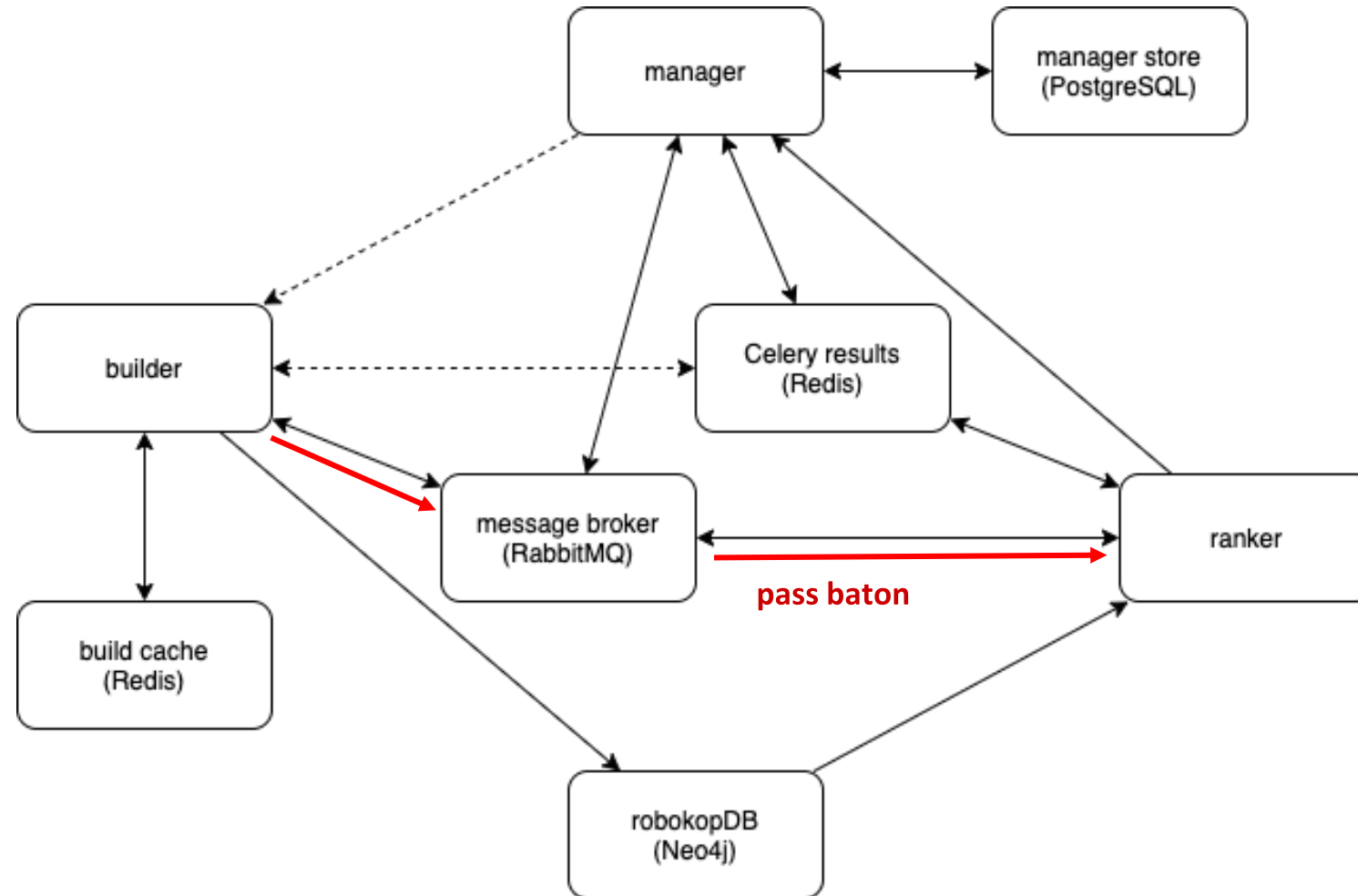  - RabbitMQ

robokop.renci.org

# ROBOKOP Containers

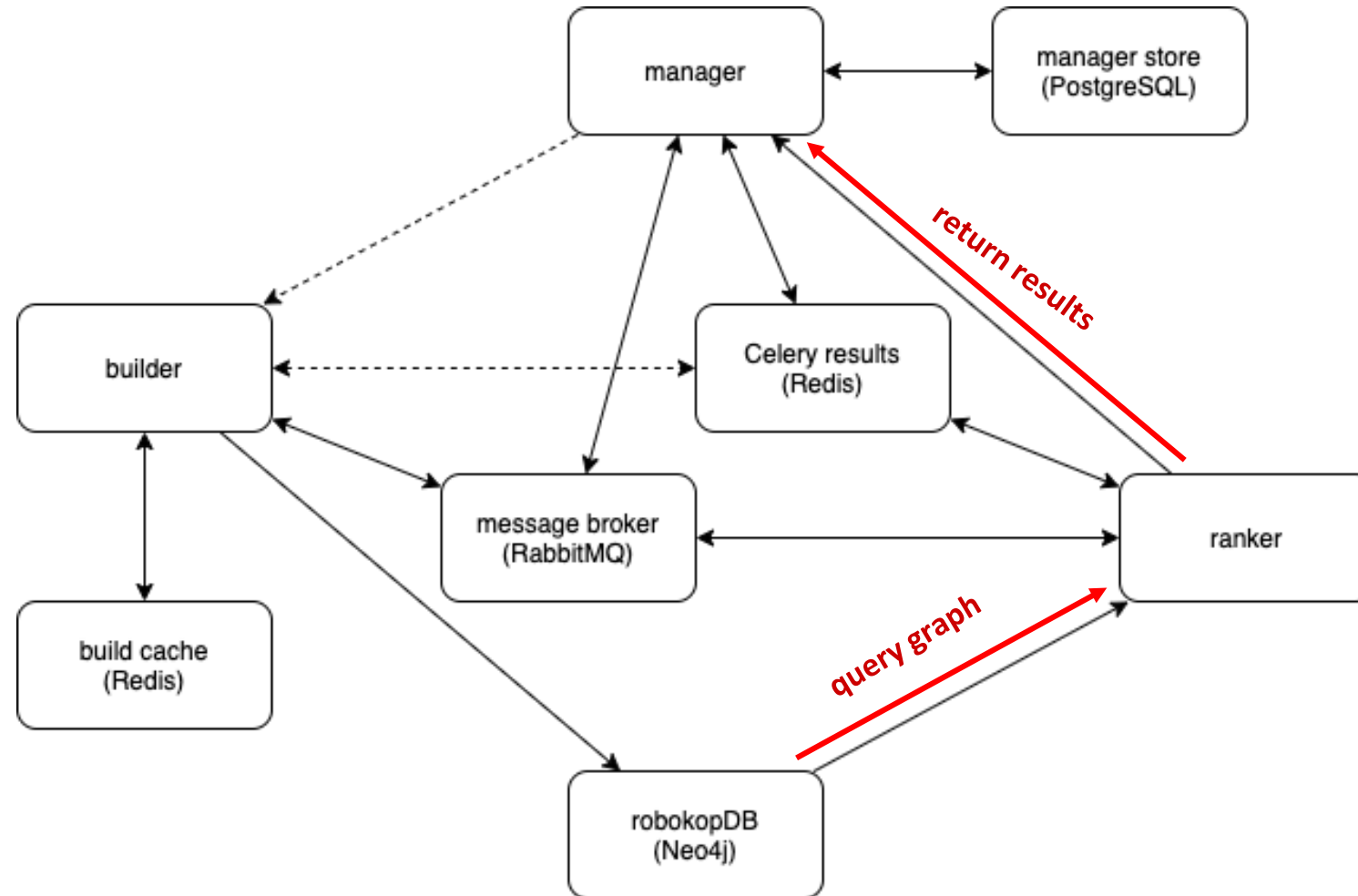# ROBOKOP Containers

# ROBOKOP Containers

# ROBOKOP Containers

# ROBOKOP Containers



robokop.renci.org

# ROBOKOP Containers



robokop.renci.org

# ROBOKOP Containers

# Shortcomings

- Brittle. If any of the following gets overworked, ROBOKOP becomes unusable:
  - robokopDB (Neo4j) - heap memory, disk space, ???
  - manager (Python) - SQLAlchemy, Celery workers
  - ranker (Python) - Celery workers
  - message broker (RabbitMQ) - *this has not happened yet!*
  - caches (Redis) - *this has not happened yet!*
  - manager store (PostgreSQL) - *this has not happened yet!*
- Opaque. Logging multiple processes, multiple containers is a mess
  - permissions for mapped Docker volumes on Linux are tricky

robokop.renci.org

# Future Work: ARAGORN

- Translator is moving from the 3-year "feasibility" phase into the 5-year "development" phase
- Greater focus on distributed system - want to avoid monolithic KG
- Want end-user interface to be responsive: seconds, not hours…
- No need to return all answers quickly, just the best ones…

- Autonomous Relay Agent for Generation of Ranked Networks (ARAGORN)
  - a.k.a Strider
  - solution: asynchronous answering

# ARAGORN

- "builder" and "ranker" get combined into a single component that builds and ranks iteratively to generate the best results first
- Plus novel algorithms for augmenting queries and "coalescing" results

- ROBOKOP was made efficient with caching and multiple processes
- ARAGORN will be efficient thanks to prioritization and asynchrony
  - We are largely I/O-bound, communicating with "third-party" Translator services

# ARAGORN Containers

# ARAGORN Tooling

- WSGI → ASGI
- Gunicorn → uvicorn
- Flask → Sanic/Starlette/FastAPI
  - asyncio, aioredis, aiormq, httpx, uvloop


- RabbitMQ + Redis for coordinating "threads"
- Neo4j for storing, querying biomedical data

# Deployment

- ROBOKOP is essentially stateful
  - Relies on a single monolithic KG
  - Relies on long-term storage for questions/answers in PostgreSQL
  - Each component scales independently
- ARAGORN is at least *less* stateful
  - Relies on short-term storage using job-specific caches
  - Easier to scale horizontally - can replicate entire ARAGORN service and parallelize behind load balancer
    - e.g. Kubernetes

# Gamma Team

Jim Balhoff
Mohit Bansal[†]
Chris Bizon (co-PI)
Richard Bruskiewich[†]
Stephen Capuzzi
Steven Cox
Miles Crosskey[†]
Karamarie Fecho
Sudeep Mandal
Kenneth Morton
Eugene Muratov
Kent Shefchek
Andrew Thieme
Alexander Tropsha (PI)
Patrick Wang

# Green Team

Stanley Ahalt (PI)
Sarav Arunachalam
Stephen Appold
Jim Balhoff
Kira Bradford[†]
James Champion
Steve Cox
Karamarie Fecho
Karl Gustafson
Adel Hanna[†]
Ray Idaszak[†]
Ann Moss Joyner[†]
Ashok Krishnamurthy
Benjamin Marsh[†]
Asia Mieczkowska[†]
Allan Parnell[†]
Dave Peden
Emily Pfaff
Kimberly Robasky[†]
Charles Schmitt
Michael Stealey[†]
Lisa Stillwell
Alexander Tropsha (co-PI)
Alejandro (Alex) Valencia
Hao Xu
Dongmei Yang[†]
Hong Yi[†]

[†]Former team member

# Funding

# References and Resources

Green/Gamma Translator Documentation Website
- ROBOKOP page
- ROBOKOP APIs

Austin CP, Colvis CM, Southall NT. Deconstructing the translational tower of babel. *Clin Transl Sci* 2019;12(2):85. doi:10.1111/cts.12595.

Bizon C, Cox S, Balhoff J, Kebede Y, Wang P, Morton K, Fecho K, Tropsha A. ROBOKOP KG and KGB: Integrated Knowledge Graphs from Federated Sources. *J Chem Inf Model* 2019 Dec 23;59(12):4968–4973. doi: 10.1021/acs.jcim.9b00683.

Morton K, Wang P, Bizon C, Cox S, Balhoff J, Kebede Y, Fecho K, Tropsha A. ROBOKOP: An Abstraction Layer and User Interface for Knowledge Graphs to Support Question Answering. *Bioinformatics* 2019;pii:btz604. doi: 10.1093/bioinformatics/btz604. [Epub ahead of print]

The Biomedical Data Translator Consortium. The Biomedical Data Translator program: conception, culture, and community. *Clin Transl Sci* 2019;12(2):86–90. doi:10.1111/cts.12592.

The Biomedical Data Translator Consortium. Toward a universal biomedical data translator. *Clin Transl Sci* 2019;12(2):91–94. doi:10.1111/cts.12591.

Contact information:
- kfecho@copperlineprofessionalsolutions.com
- patrick@covar.com