

We're trying to predict FUTURE discoveries by reading the ENTIRE history of astronomy papers! (No pressure...)

PREDICTING NEW CONCEPT-OBJECT ASSOCIATIONS IN ASTRONOMY BY MINING THE LITERATURE

Jinchu Li
School of Computer Science,
Georgia Institute of Technology
Atlanta, GA 30332, USA
jinchu.li@gatech.edu

Yuan-Sen Ting
Department of Astronomy
The Ohio State University
Columbus, OH 43210, USA
ting.748osu.edu



Literature-based Crystal Ball

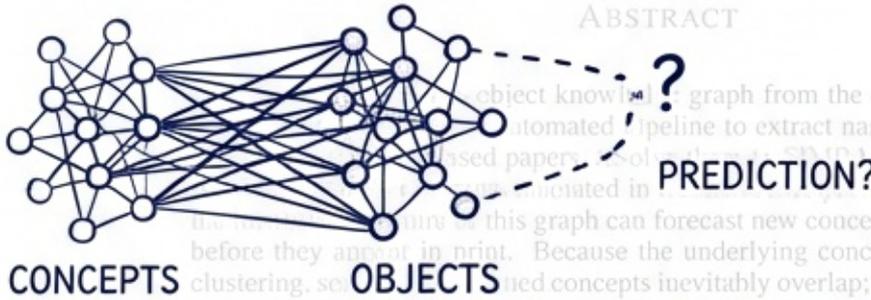
Alberto Accomazzi

Tirthankar Ghosal

TL;DR SUMMARY

We built a GIANT "knowledge graph" linking astronomical OBJECTS (like stars, galaxies) to scientific CONCEPTS from tons of old papers.

Then we asked: Can the past connections help us guess NEW connections before they get published?



Spoilers: YES!
Our method is way better at guessing than just looking at what's similar.



This could save astronomers a TON of time & money by telling them where to look first!

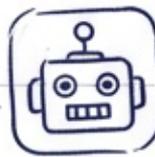
WHY BOTHER?

- difficulty of navigating a rapidly growing literature
- undiscovered public knowledge

Because science is **DROWNING** in papers! Too many for any human to read.



The answers might be hidden in plain sight, connecting two things nobody thought to put together!

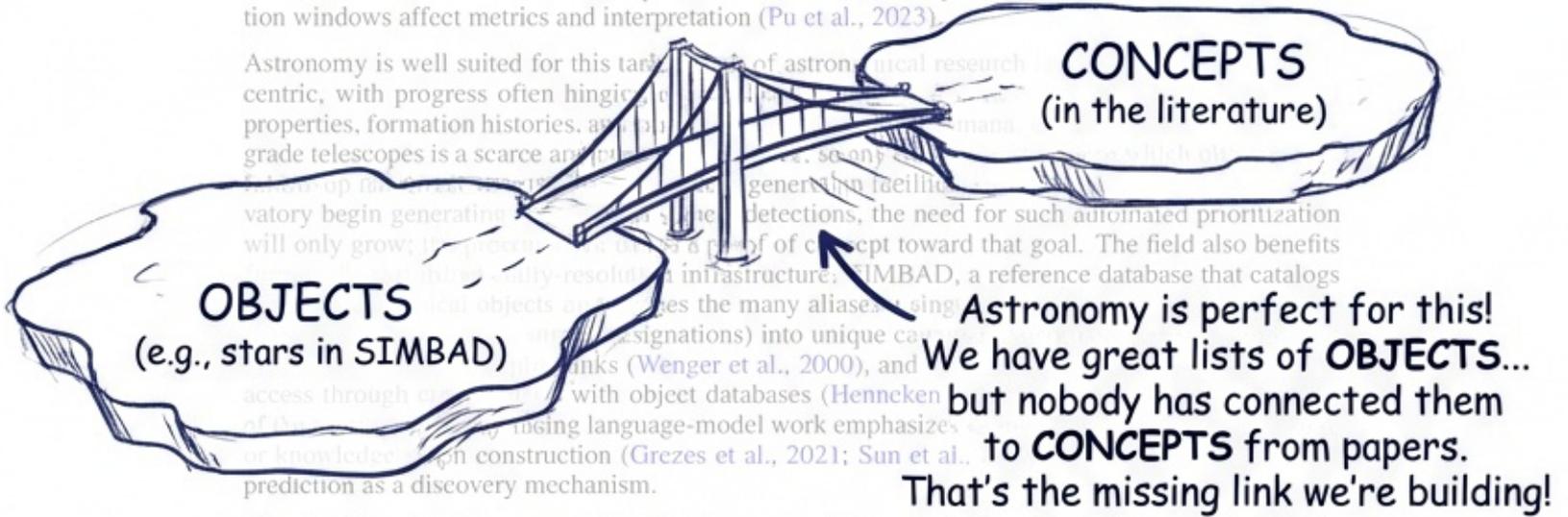


So, what changed? **Better AI for reading!** Now computers can actually "understand" scientific text and build these complex knowledge graphs. (Before, it was just too hard!)

What has changed recently is the feasibility of building richer representations of scientific knowledge from text. Pretrained and large language models have improved extraction and representation for scientific corpora, enabling scalable pipelines for entity recognition, normalization, and relation mining (Devlin et al., 2019; Beltagy et al., 2019; Brown et al., 2020), and have been used directly for knowledge graph completion and inductive link prediction (Yao et al., 2019; Daza et al., 2021). These advances shift the bottleneck from whether we can construct literature-derived graphs to how we evaluate and use them for discovery.

Link forecasting on literature-derived graphs has accordingly become an active research direction. SemNet constructs evolving concept networks from scientific text and tests whether historical snapshots predict future connections (Krenn & Zeilinger, 2020); Science4Cast standardizes the protocol of training up to a temporal cutoff and predicting which links appear afterward (Aghajohari et al., 2021). More recently, Impact4Cast scales forecasting to tens of millions of papers (Gu & Krenn, 2024), while biomedical studies adapt LBD to time-sliced link prediction and examine how evaluation windows affect metrics and interpretation (Pu et al., 2023).

Astronomy is well suited for this task. The field of astronomical research is concept-centric, with progress often hinging on the discovery of new objects and their properties, formation histories, and interactions. The discovery of the first space-grade telescopes is a scarce and expensive endeavor, and the generation of high-quality data only grows with time. The need for such automated prioritization is a proof of concept toward that goal. The field also benefits from the construction of infrastructure: SIMBAD, a reference database that catalogs objects (Wenger et al., 2000), and ADS, a database that links papers to objects (Henneken et al., 2019). Building such a graph requires bridging the gap between object databases (Henneken et al., 2019) and concept extraction and resolution across the full literature.



Despite these databases, no systematically constructed concept-object knowledge graph exists in astronomy, to our knowledge. SIMBAD catalogs objects but does not map them to the scientific concepts they relate to; ADS links papers to objects but does not map them to the scientific concepts they relate to. Building such a graph requires bridging the gap between object databases (Henneken et al., 2019) and concept extraction and resolution across the full literature.

PIONEERING WORK

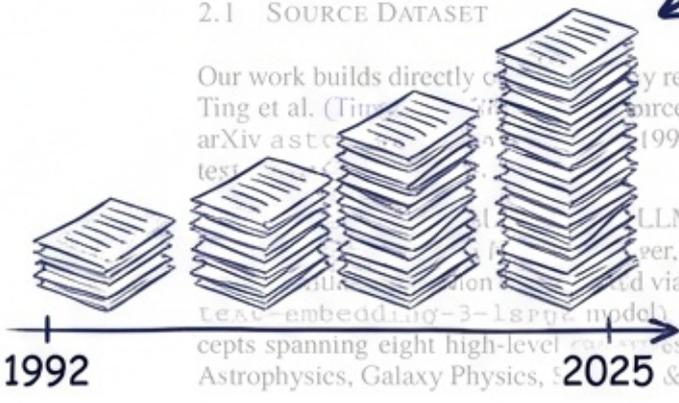
Our work addresses this gap. We build an end-to-end pipeline that maps object mentions to SIMBAD identifiers and links them to concepts across the full literature. As a proof of concept, we evaluate this pipeline on a task: given all associations observed up to a cutoff year, we test whether a collaborative-filtering model can predict which objects will newly appear in the literature with a given concept. Figure 1 provides an overview.

2 DATASET

2.1 SOURCE DATASET

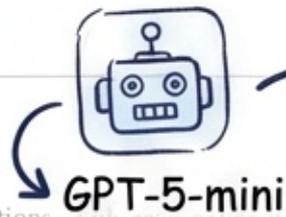
Our work builds directly on the source corpus comprises 408,590 astrophysics papers from the arXiv astrophysics preprint server, released in 1992 through July 2025, each converted to machine-readable text. We use an LLM-based extraction pipeline to identify key scientific and concepts spanning eight high-level categories (e.g., Cosmology & Nongalactic Physics, High Energy Astrophysics, Galaxy Physics, & AI). The

To do this, we used A LOT of data. **Over 400,000 astronomy papers!** (from arXiv, spanning nearly 30 years)



We let an LLM loose on all these papers to automatically find the key concepts. (Because who has time to read 400k papers manually?!)





Okay, we have the CONCEPTS.
Now we need the OBJECTS.
We ask the AI (GPT-5-mini) to read
the papers and find them for us.

paper-concept associations, with each concept accompanied by a descriptive text definition and a fixed embedding vector.

We treat this concept vocabulary, its associated embeddings, and the paper-concept links as fixed inputs, and focus on augmenting the dataset with a new object-extraction procedure. It also tells us how the object was used (was it the main subject? just a comparison?).

2.2 OBJECT EXTRACTION AND SIMBAD RESOLUTION

For each paper, we prompt GPT-5-mini with the title, abstract, and full OCR text to extract candidate astronomical object mentions intended for SIMBAD-style resolution. The extraction procedure is summarized in Appendix A; the full prompt and output schema are available in the Data and Code Availability section. For each extracted object, the model returns: (i) a single designation string, (ii) a semantic role describing how the object is used in the paper (e.g., primary subject, sample member, counterpart/host, calibration/reference), (iii) a study mode (new observations, archival/reanalysis, catalog compilation, theory/simulation, or incidental mention), and (iv) a short evidence span with its source (title/abstract/body).

Raw AI output

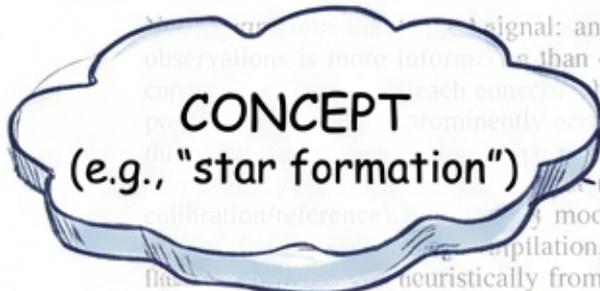


But AI can make mistakes.
So we pass its answers through a "truth filter" (SIMBAD database) to make sure the objects are real and we have their correct names.

the remaining object strings are filtered to remove duplicates. We validate against the SIMBAD database (a fail validation constraint) and deduplicate the remaining object strings against the canonical identifiers (Wenger et al. 2000); names that mention instances removed. Across 501 papers, the remaining 1,621,446 mentions remain (376,672 unique names), of which 163,775 unique names (43.48%) successfully resolve, corresponding to 1,150,553 surviving mention instances (66.06% of raw mentions). After alias merging, we retain 100,560 unique SIMBAD objects across the corpus.

2.3 CONCEPT-OBJECT GRAPH CONSTRUCTION

We combine the paper-concept associations from the source dataset with the paper-object associations extracted and resolved in the previous step to form a concept-object bipartite graph. An edge between concept c and object o exists if and only if at least one paper mentions o in connection with c . Each mention inherits the publication year of its source paper, enabling temporal analysis.



Strong link: Object was the MAIN subject of the paper.



Weak link: Object was just mentioned in passing.

signal: an object that is the primary observation is more informative than one mentioned only in passing. Each concept-object edge a weight that reflects how many papers prominently mention figures in its paper. As described in Section 2.2, role captures the semantic role and a study mode. Role captures the semantic role (e.g., primary scientific target vs. mere calibration/reference) and study mode captures the type of analysis (e.g., archival/reanalysis, theory-only, or catalog compilation, theory-only, or incidental mention). We filter mentions heuristically from domain information—for example, mentions via new observations should receive more weight, compared to a mere mention in a theoretical context.

Formally, the weight of a concept-object edge is computed by aggregating per-mention weights and applying a log transform:

where $\mathcal{M}(c, o)$ is the set of paper-level mentions linking concept c to object o . Each mention's strength is the product of a role weight and a study-mode weight:

$$a(m) = \rho_{r(m)} \gamma_{o(m)}, \quad (2)$$

where $r(m)$ and $o(m)$ are the role and study mode of mention m , and ρ and γ are fixed lookup weights ($\rho \in [0.25, 3.0]$, $\gamma \in [0, 1.25]$ for roles and study modes).

This scary-looking math is just a way to calculate that link strength based on how important the object was in the paper.

PUTTING IT ALL TOGETHER: THE COMPLETE PIPELINE!

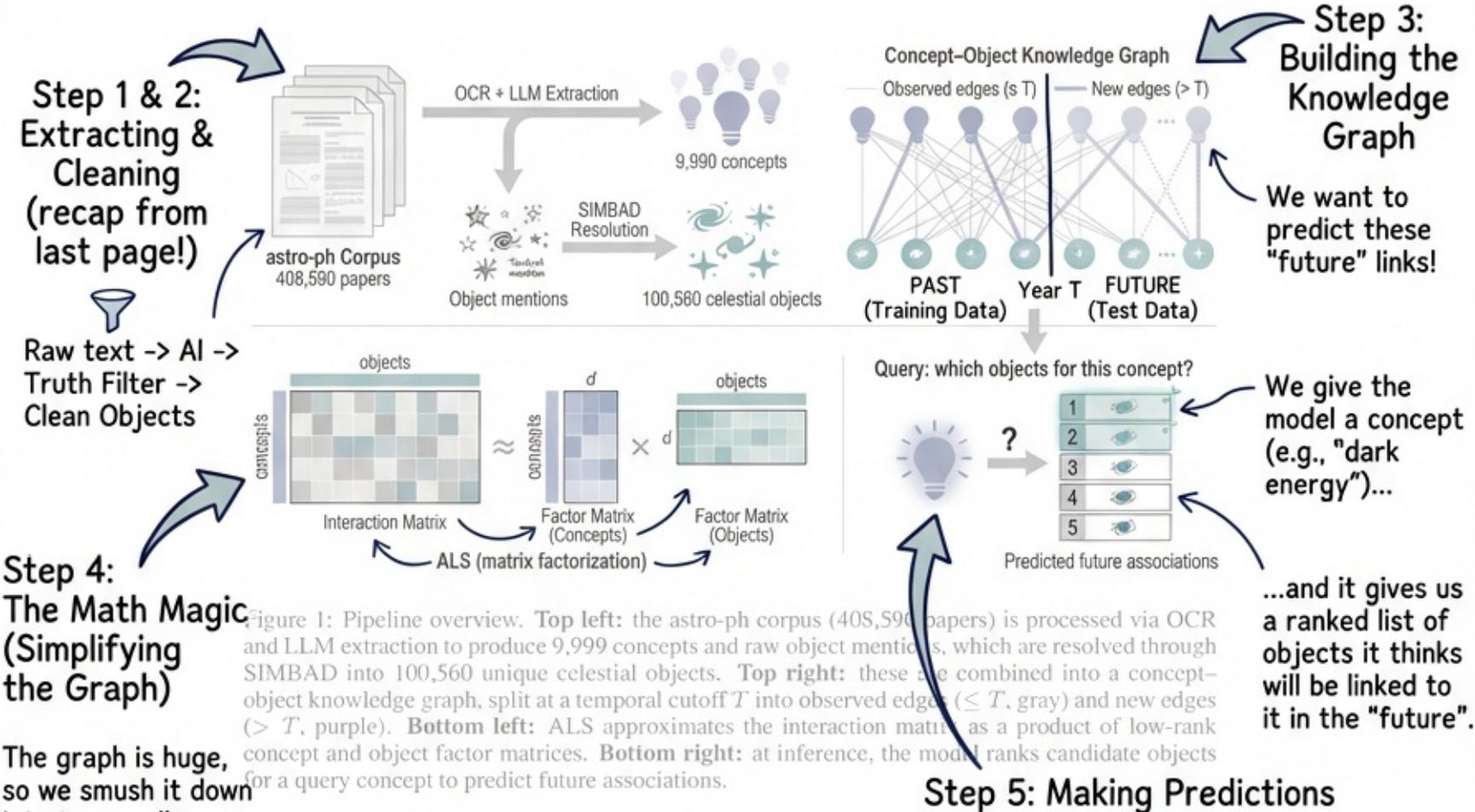


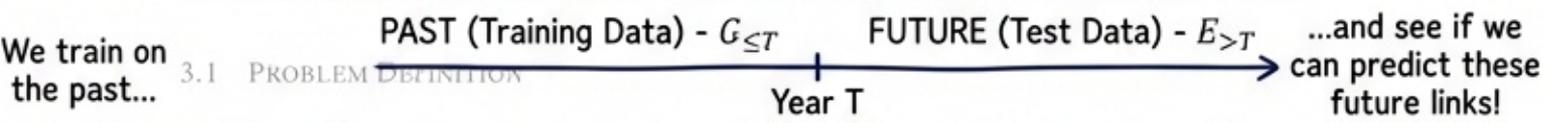
Figure 1: Pipeline overview. Top left: the astro-ph corpus (408,590 papers) is processed via OCR and LLM extraction to produce 9,999 concepts and raw object mentions, which are resolved through SIMBAD into 100,560 unique celestial objects. Top right: these are combined into a concept-object knowledge graph, split at a temporal cutoff T into observed edges ($\leq T$, gray) and new edges ($> T$, purple). Bottom left: ALS approximates the interaction matrix as a product of low-rank concept and object factor matrices. Bottom right: at inference, the model ranks candidate objects for a query concept to predict future associations.

The log transform compresses the dynamic range so that edges supported by many weak mentions do not dominate edges supported by a few strong ones. We additionally record a timestamp for each concept-object pair, defined as the earliest publication year among its supporting papers; this timestamp is what enables the temporal train/test split used in our evaluation (Section 3.1). Representative edges from the concept-object graph for two illustrative concepts, including aggregated edge weights and first-appearance years, are shown in Appendix C.

3 METHODS

How do we know this works? We test its ability to "predict the future"!

The central question is whether the concept-object graph contains enough structure to anticipate future associations. We adopt a temporal holdout protocol: freeze the graph at a cutoff year T , train on all associations up to that point, and evaluate predictions of associations that first appear after T .



Let C denote the set of concepts and O the set of resolved astrophysical objects. From the graph $G = (C, O, E)$ constructed in Section 2, an edge $(c, o) \in E$ carries a positive weight $w_{c,o} > 0$ and a first-appearance year $y_{c,o}$.

Given a cutoff year T , we split the graph temporally: the training graph $G_{\leq T}$ contains only edges whose first-appearance year is at most T , with edge set $E_{\leq T}$, while the held-out set $E_{>T}$ contains edges appearing strictly after T . For each concept c , the task is to rank candidate objects $O_c^{\text{cand}} = \{o \in O \mid (c, o) \notin E_{\leq T}\}$ —that is, objects not yet associated with c —so that the ones in $E_{>T}$ are ranked highly. For example, if the concept "High-Redshift Quasars" has been linked to a set of known quasars before T , the model should rank the quasars that will first appear in connection with this concept after T above the remaining candidates.

e.g., The model should know beforehand which quasars will become important for this concept.

THE CORE ENGINE: MAKING RECOMMENDATIONS



3.2 MODEL AND BASELINES

The ranking problem defined above can be viewed as a recommendation task: given the historical concept-object associations, recommend which objects a concept is likely to be associated with next.

We treat this as an implicit-feedback problem, meaning we observe which concepts co-occurred (positive signal) but never observe explicit “negative” pairs—going to an object may simply reflect an association that has not yet appeared in past papers. To work with standard matrix methods, we represent the graph as a matrix $W \in \mathbb{R}_{\geq 0}^{(|C|) \times (|O|)}$, where $W_{c,o} = w_{c,o}$ for observed edges and $W_{c,o} = 0$ for predicted edges. The goal is to learn a scoring function that assigns high scores to objects that become nonzero in the future.

Think of this like Netflix recommending your next binge-watch. We have patterns of past connections (positive signal) and want to predict *new hits* that haven't happened yet.

Matrix Factorization. Our primary model is implicit Alternating Least Squares (ALS) (Hu et al., 2008), a matrix factorization method. The core idea is that the interaction matrix W , despite being very large (9,999 concepts \times 100,300 objects), has low-rank structure: the patterns of which concepts associate with which objects are governed by a much smaller number of latent factors (e.g., shared physical properties, observational techniques, or research themes). We approximate W by representing each concept c as a vector $p_c \in \mathbb{R}^d$ and each object o as a vector $q_o \in \mathbb{R}^d$ (where $d \ll |C|, |O|$ is the latent dimension), such that their dot product $p_c^\top q_o$ estimates the strength of their association. Because concepts that share similar objects are pushed toward similar vectors during training, the model can generalize to predict associations for concept-object pairs not seen in the training data.

Concretely, the model minimizes the squared error between observed edges and predicted edges with ℓ_2 regularization:

THE MATH BEHIND THE MAGIC (The Objective Function)

$$\min_{\{p_c\}, \{q_o\}} \sum_{c \in C} \sum_{o \in O} (1 + \alpha w_{c,o}) (\mathbb{I}[w_{c,o} > 0] - p_c^\top q_o)^2 + \lambda \left(\sum_c \|p_c\|_2^2 + \sum_o \|q_o\|_2^2 \right), \quad (3)$$

where $\mathbb{I}[w_{c,o} > 0]$ is a binary indicator of whether an association has been observed, $\alpha = 10$ is a hyperparameter controlling how much the weight amplifies confidence, and $\lambda = 0.05$ scales the regularization strength. The model is trained on a baseline of concept-aware KNN and an inference-time smoothing procedure. The predicted relevance score for a concept-object pair is then $s_{\text{ALS}}(c, o) = p_c^\top q_o$.

Try to match the links we've actually seen (weighted by confidence)...but keep the model simple so it doesn't just memorize the data (“overfitting”).

Concept-embedding similarity weights. Because the 9,999 concepts in the source dataset are semantically related, concept embeddings inevitably overlap, fragmenting the space of concept embeddings. We use concept-aware KNN and an inference-time smoothing procedure to account for this overlap; these weights serve as a concept-aware KNN baseline and an inference-time smoothing procedure. For each concept c and define nonnegative neighbor weights by clipping negative cosines and ℓ_1 -normalizing:

$$S_{c,c'} = \frac{\max(\cos(c_c, c_{c'}), 0)}{\sum_{c'' \in \mathcal{N}_k(c)} \max(\cos(c_c, c_{c''}), 0)}, \quad (4)$$

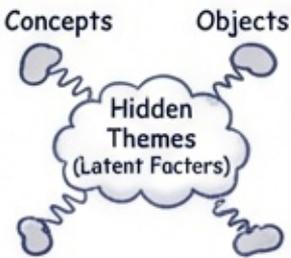
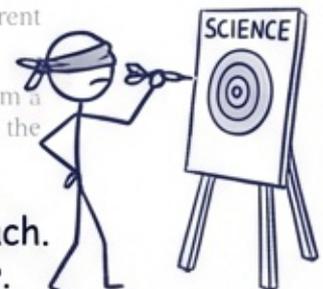
for $c' \in \mathcal{N}_k(c)$, and $S_{c,c'} = 0$ otherwise. If all similarities in $\mathcal{N}_k(c)$ are non-positive, we fall back to uniform weights over $\mathcal{N}_k(c)$.

THE COMPETITION (Things we have to beat)

Baselines. We compare ALS against several non-parametric baselines, each testing a different hypothesis about what drives future associations.

Random ranks candidate objects uniformly at random, assigning each object a score drawn from a uniform distribution: $s_{\text{rand}}(o) \sim \text{Uniform}(0, 1)$. This establishes a lower bound and verifies that the evaluation protocol behaves as expected.

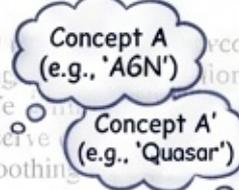
1. **RANDOM GUESSING:** The “monkey throwing darts” approach. Sets the absolute bottom bar.



Back to the “Math Magic” (ALS) from the last page! The assumption: deep down, there are hidden structures connecting things.

HANDLING SIMILAR IDEAS: Like so many other things, sometimes different words mean almost the same thing. This helps the model share info between them.

HANDLING SIMILAR IDEAS: Sometimes different words mean almost the same thing. This helps the model share info between them.



The Contenders & The Rules of the Game

The Lineup of Me're comparing

How we keep score

The 'monkey' baseline (see prev. page)

The 'Bandwagon' approach: just recommend what's already famous.

The 'Trendy' approach: what's famous lately.

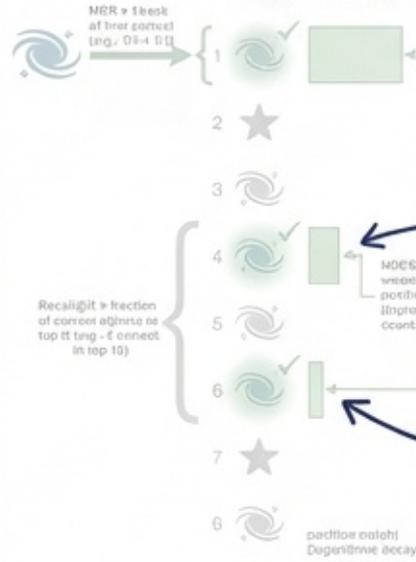
The 'Copycat' approaches: find similar concepts and borrow their associations.

Our 'Matrix Factorization' star (from prev. page).

Baselines & Model



Evaluation Metrics



Did you get the right answer *quickly*? (Higher is better).

Did you find *all* the right answers in your top K picks? (Coverage).

The gold standard: Cares about finding the right answers AND putting them at the top.

Figure 2: Visual overview of baselines and evaluation metrics. **Left:** how each method ranks candidate objects for a query concept—from top to bottom: Random (shuffled ordering), Popularity (global frequency), RecentPopularity (time-windowed frequency), ConceptKNN-AA (Adamic-Adar neighbor aggregation), ConceptKNN-TextEmb (embedding-based neighbor aggregation), and ALS (dot product of learned latent factors). **Right:** how the three evaluation metrics score a ranked list; green checkmarks denote correct held-out associations. MRR rewards placing the first correct object high. Recall@ K measures coverage in the top K , and NDCG@100 assigns position-discounted credit.



Popularity tests whether globally prominent objects are more likely to appear in future associations regardless of the query concept, scoring each object by its total training edge weight $s_{pop}(o) = \sum_{c \in \mathcal{C}} w_{c,o}$. The resulting ranking is concept-agnostic.

RecentPopularity ($\Delta \in \{3, 5\}$ years) restricts the sum to papers published in $(T-\Delta, T]$, testing whether short-term trends carry more signal than lifetime popularity. Like Popularity, it is concept-agnostic.

ConceptKNN family. We define a concept-aware non-parametric family of baselines that borrow evidence from similar concepts. For a query concept c , we form a neighbor set $\mathcal{N}_b(c)$ under a chosen concept-concept similarity and score candidate objects by

$$s_{knn}(c, o) = \sum_{c' \in \mathcal{N}_b(c)} \tilde{s}(c, c') w_{c', o}, \quad (5)$$

where $\tilde{s}(c, c')$ are nonnegative neighbor weights normalized over $\mathcal{N}_b(c)$. We treat k as a hyperparameter and report results across multiple values. We use two instantiations:

(i) **ConceptKNN-AA (graph-based).** We set similarity using the Adamic-Adar index (Adamic & Adar, 2005):

$$s_{AA}(c, c') = \sum_{o \in \mathcal{O}(c) \cap \mathcal{O}(c')} \frac{1}{\log |\mathcal{C}(o)|}, \quad (6)$$

and define $\tilde{s}(c, c')$ by normalizing $s_{AA}(c, c')$ over $c' \in \mathcal{N}_b^k(c)$.

(ii) **ConceptKNN.Emb (text-embedding-based).** We set $\tilde{s}(c, c') = S_{c,c'}$, the concept-embedding similarity weights defined in equation 4. This baseline scores objects *entirely* by aggregating the associations of semantically similar concepts.

Training and inference. For each cutoff T , we train all methods on the time-filtered graph $G_{\leq T}$, keeping the concept and object vocabularies fixed across cutoffs. The ALS model is fit on the

Lazy! Gives the same recommendations to everyone, regardless of the concept (concept-agnostic).

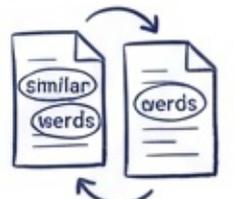


Smarter! Tailors recommendations based on the specific concept (concept-aware).



Uses shared connections in the data.

Uses similar meanings from text descriptions.



A quick note on fair play: we train on past data and test on future data. Standard ML practice.



Final "Rules of the Game" & The Official Scoring Rulebook

A few more →
"ground rules" to
keep the fight
fair...



Fine-tuning our
ALS model's "engine"
before the race.
We used a small
grid search to find
the best settings.

corresponding interaction matrix to produce latent vectors for all concepts and objects; baselines compute their scores directly from training-graph statistics. At inference time, for each concept c we rank all objects in the candidate set O_c^{cand} by their predicted score $s(c, o)$. The same candidate masking—excluding objects already associated with c in $G_{<T}$ —is applied to all methods, ensuring a fair comparison. We select ALS hyperparameters (d, α, λ) via a small grid search on a temporally valid development setting (using cutoff $T=2021$ for selection), and then fix $d=125$, $\alpha=10$, and $\lambda=0.05$ for all reported cutoffs. For ConceptKNN baselines, we report results across $k \in \{5, 10, 25, 50, 100, 150, 200\}$. Additionally, the concept-embedding weights enable smoothing any method's predictions at inference time. Given a base score $s(c, o)$, we blend it with neighbor-aggregated scores:

$$s_{smooth}(c, o) = (1 - \beta) s(c, o) + \beta \sum_{c' \in \mathcal{N}_k(c)} S_{o, c'} s(c', o), \quad (7)$$



Giving the "Copycat" (ConceptKNN) a helping hand: This "smoothing" blends its score with its neighbors, making it a stronger competitor.

where β controls the mixing strength. This preserves each method's learned structure while compensating for concept-boundary artifacts. We apply smoothing uniformly to all methods, tuning (k, β) at $T=2021$ and fixing $k=100$, $\beta=0.5$ for all reported experiments. Results without smoothing are in Appendix D.

3.3 EVALUATION PROTOCOL

Evaluation is restricted to concepts with at least 10 training associations and at least one held-out test edge in $E_{>T}$. The first condition ensures that concepts have enough history for the model to learn from; the second ensures a meaningful prediction target.

We assess ranking quality with four standard metrics, where C_{eval} is the $E_{>T}(c)$ the held-out objects for concept c :

Mean Reciprocal Rank (MRR): the reciprocal rank of the first correct held-out object, measuring how quickly a correct result appears:

$$MRR = \frac{1}{|C_{eval}|} \sum_{c \in C_{eval}} \frac{1}{\min_{o \in E_{>T}(c)} \text{rank}_v(o)}. \quad (8)$$



Remember: Speed is key! 1st place is WAY better than 2nd.

Recall@K (with $K \in \{10, 100\}$): the fraction of held-out objects appearing in the top K predictions:

$$\text{Recall@K} = \frac{1}{|C_{eval}|} \sum_{c \in C_{eval}} \frac{|\{o \in E_{>T}(c) : \text{rank}_v(o) \leq K\}|}{|E_{>T}(c)|}. \quad (9)$$



Did your net catch the right fish? (How many good answers are in your top K?)

NDCG@100: a ranking-quality measure that awards more credit for correct objects placed near the top (Järvelin & Kekäläinen, 2002):

$$NDCG@100 = \frac{1}{|C_{eval}|} \sum_{c \in C_{eval}} \frac{DCG@100(c)}{IDCG@100(c)}, \quad DCG@100(c) = \sum_{i=1}^{100} \frac{\mathbb{I}[\mathcal{R}_o(i) \in E_{>T}(c)]}{\log_3(i+1)}, \quad (10)$$

where

$$IDCG@100(c) = \sum_{i=1}^{\min(|E_{>T}(c)|, 100)} \frac{1}{\log_2(i+1)} \quad (11)$$

is the ideal DCG obtained by ranking all held-out objects first. All metrics are macro-averaged over concepts; no test-edge information is used during training.

Are the BEST answers at the very TOP? Ranking order *really* matters here.

4 RESULTS

We evaluate concept-object edge prediction across four cutoff years $T \in \{2017, 2019, 2021, 2023\}$. Depending on the cutoff, the training graph contains 2.11M–2.90M concept-object edges, and the held-out set contains 0.26M–1.05M edges. Evaluation is restricted to 7.2k–7.7k eligible concepts per cutoff (each with at least 10 training associations and at least one held-out test association), with a fixed object vocabulary of 100,560 SIMBAD-resolved objects.



Setting the stage for the MAIN EVENT! We're testing on a MASSIVE scale across different years to be thorough.

Running the race at different points in time!

The Official Scoring Rulebook (The Math-y version of the previous page!)



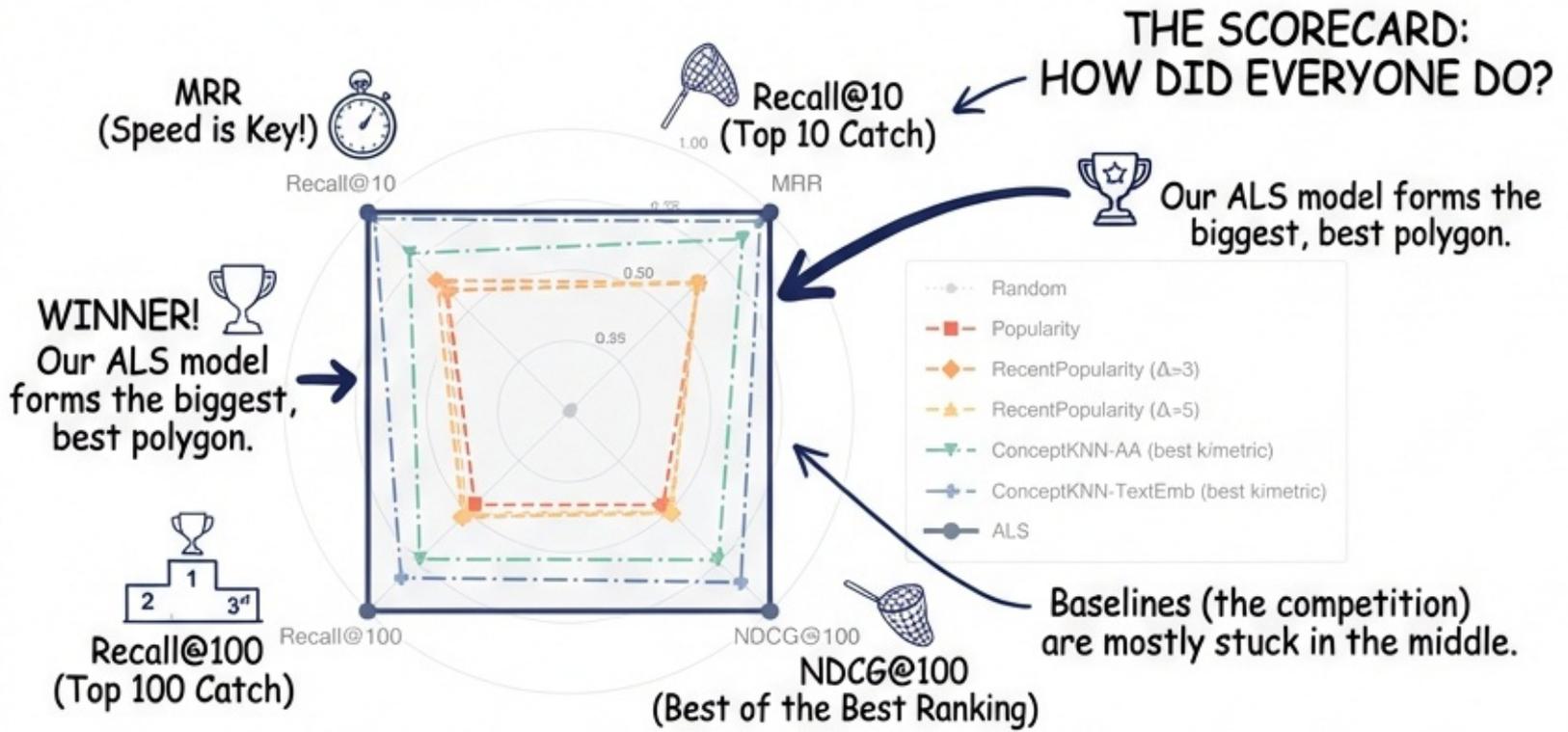


Figure 3: Radar plot comparing all methods on the physical concept subset with concept smoothing. The four axes correspond to MRR, Recall@10, Recall@100, and NDCG@100, each normalized so the best method equals 1.0. Methods shown: Random (gray dotted), Popularity (red dashed), RecentPopularity with $\Delta \in \{3, 5\}$ (orange/yellow dashed), ConceptKNN-AA (teal dash-dot), ConceptKNN-TextEmb (blue dash-dot), and ALS (dark solid). ALS forms the outermost polygon, leading on all four metrics.

THE DETAILED SCORE SHEET
 (The Nitty Gritty)

(a) Absolute performance (with inference-time concept smoothing)

Method	MRR	Recall@10	Recall@100	NDCG@100
Random	0.0061±0.0020	0.0024±0.0001	0.0010±0.0001	0.0011±0.0003
Popularity	0.2031	0.0278	0.0817	0.0673
RecentPopularity ($\Delta=3$)	0.2031	0.0303	0.0827	0.0734
RecentPopularity ($\Delta=5$)	0.2024	0.0282	0.0902	0.0719
ConceptKNN AA (best k per metric)	0.2724	0.0365	0.1294	0.1070
ConceptKNN-TextEmb (best k per metric)	0.3040	0.0445	0.1457	0.1230
ALS	0.3150±0.0010	0.0460±0.0001	0.1746±0.0002	0.1436±0.0001

Our Champion (ALS) dominates on EVERY single metric!

(b) Relative improvement of ALS over baselines (%)

Baseline	MRR	Recall@10	Recall@100	NDCG@100
Random	5054.48	20694.20	18176.33	12756.58
Popularity	55.13	65.43	113.75	112.66
RecentPopularity ($\Delta=3$)	55.11	51.81	88.31	95.63
RecentPopularity ($\Delta=5$)	55.62	62.90	93.60	99.55
ConceptKNN-AA (best k per metric)	15.62	26.05	34.94	34.20
ConceptKNN-TextEmb (best k per metric)	5.01	3.39	19.80	16.78

The "dart-throwing monkey" (Random) is predictably terrible.

This table shows the improvement of ALS over the others. It's often HUNDREDS or THOUSANDS of percent better than the basic baselines! Even the "tougher competitors" (ConceptKNN) get left behind by 15-20%.

Table 1: Link prediction performance on the physical concept subset with concept smoothing, averaged across four cutoffs. ALS and Random are stochastic and reported as mean±std over seeds; all other baselines are deterministic. KNN baselines: best k per metric. Panel (b): relative improvement of ALS over each baseline (%).

Evaluation concept subset. Not all concept categories are equally informative for forecasting astrophysical associations. Concepts tied to instrument usage or broadly applicable methodologies yield associations driven more by scheduling or cross-field applicability than by astrophysical structure, diluting evaluation. We therefore report headline results on a *Physical* subset that excludes Statistics & AI, Numerical Simulation, and Instrumental Design concepts (details in Appendix E).

Focusing on the REAL science. We excluded concepts like

8 "telescope" or "method" to see how well we predict actual astrophysical ideas. This makes the test harder and more meaningful!

THE TEST OF TIME:

Does it hold up year after year?

Method	2017	2019	2021	2023
	Random	0.0010±0.0002	0.0009±0.0002	0.0008±0.0003
Popularity	0.0764	0.0812	0.0865	0.0826
RecentPopularity ($\Delta=3$)	0.0841	0.0846	0.1001	0.1021
RecentPopularity ($\Delta=5$)	0.0816	0.0845	0.0964	0.0982
ConceptKNN-AA (best k)	0.1201	0.1317	0.1412	0.1430
ConceptKNN-TextEmb (best k)	0.1324	0.1429	0.1547	0.1594
ALS	0.1463±0.0004	0.1596±0.0004	0.1766±0.0005	0.1851±0.0005

(b) NDCG@100

Method	2017	2019	2021	2023
	Random	0.0015±0.0003	0.0012±0.0003	0.0009±0.0002
Popularity	0.0832	0.0744	0.0647	0.0478
RecentPopularity ($\Delta=3$)	0.0909	0.0771	0.0704	0.0352
RecentPopularity ($\Delta=5$)	0.0882	0.0775	0.0652	0.0336
ConceptKNN-AA (best k)	0.1369	0.1226	0.1063	0.0848
ConceptKNN-TextEmb (best k)	0.1365	0.1349	0.1193	0.0954
ALS	0.1617±0.0002	0.1459±0.0003	0.1306±0.0003	0.1070±0.0003

SPOILER:
YES. ALS
wins every
single year.



Table 2: Per-cutoff performance on the physical concept subset without smoothing. ALS leads at every cutoff on both metrics. ALS and Random: mean±std over seeds; other baselines are deterministic. KNN baselines: best k per cutoff.

Table 1 and Figure 3 report performance averaged across all four cutoffs with concept smoothing applied uniformly to all methods. ALS is the strongest method on all four metrics, exceeding the best text-embedding KNN baseline by 5.0% on MRR, 3.4% on Recall@10, 19.8% on Recall@100, and 16.8% on NDCG@100. The gains over popularity heuristics are substantially larger: ALS improves Recall@10 by 10% and 96% on NDCG@100. The margin is widest on long-horizon retrieval metrics (Recall@100, NDCG@100)—precisely the regime most relevant to practical triaging, where a researcher scans a moderately sized candidate list. On early-rank metrics the advantage is smaller but still positive, indicating that matrix factorization captures complementary structure beyond local neighborhood voting. Results without smoothing (Appendix D) preserve the qualitative ordering but show a narrower ALS advantage on early-rank metrics; smoothing most strongly benefits ALS, consistent with the idea that learned latent factors propagate information more effectively once cluster-boundary artifacts are mitigated.

Because aggregate improvements are largest on long-horizon metrics, we report results **without smoothing, for a cleaner comparison of base methods.** Table 2 breaks down Recall@100 and NDCG@100 by cutoff year (without smoothing, for a cleaner comparison of base methods; smoothed per-cutoff trends are in Appendix F). ALS outperforms all baselines at every cutoff on both metrics, confirming that the gains are consistent across time. RecentPopularity becomes more competitive at later cutoffs but does not close the gap; KNN baselines improve steadily with T but remain behind ALS.

Why this matters:
We're best at the "long game"—finding distant, non-obvious connections useful for "practical triaging".

We took the "smoothing" filter off here to show the RAW performance. ALS still dominates.

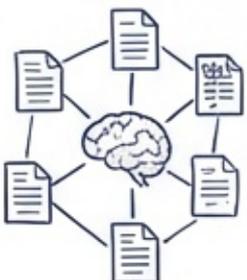
THE BOTTOM LINE: Consistent wins, especially on the hard stuff (real physics concepts & long-horizon search).

5 DISCUSSION, IMPLICATIONS

Summary of findings. With inference-time concept smoothing to mitigate the semantic overlap inherent in clustering-derived concepts, ALS consistently outperforms all baselines across every metric and temporal cutoff, with the largest gains on long-horizon retrieval measures most relevant to practical triaging. These improvements are concentrated in astrophysically meaningful concepts, indicating that the literature-derived graph captures predictive structure beyond local similarity or short-term trends.

Implications for literature-driven discovery. The big idea is that by mining the published literature can surface latent structure as SIMBAD and ADS—not by cataloging objects or their relationship between scientific ideas and the objects

BIG IDEA: The literature isn't just a pile of papers. It's a giant, evolving MAP of knowledge. Our model helps draw that map.



A FINAL WORD OF CAUTION:

Our tool finds expected connections based on the past. But the **BIGGEST** scientific leaps often come from the **UNEXPECTED!** This is a tool to help your judgment, not replace it.



building block for tools that help astro- explored objects for a given research question. **Limits of prediction and scientific judgment.** At the time of our research, we know that future concept-object associations are statistically forecastable from past literature, some of the most consequential discoveries in astronomy arise precisely from unexpected associations—objects that no one predicted would be relevant to a given research question. A triaging system that ranks objects by expected relevance risks reinforcing existing research trajectories at the expense of serendipitous findings. Balancing systematic prioritization with openness to the unexpected is not a problem that any algorithm can resolve on its own; it is ultimately a matter of scientific judgment and community practice. We view the graph presented here as a tool to inform that judgment, not to replace it.

Methodological limitations and design choices. The concept-object graph depends on OCR quality, LLM-based extraction accuracy, and SIMBAD. In addition, our evaluation targets are defined by a ground-truth notion of physical discovery, so prediction and publication dynamics.

Because the 9,999 concepts are obtained by clustering, overlap is unavoidable: closely related or near-synonymous concepts are often scientifically coherent. Our inference-time concept assignment to mitigate this, but more principled approaches—such as joint concept representations, soft assignment, or jointly learned concept embeddings—could potentially address concept overlap more directly and are a promising direction for future work.

Scope and future directions. The object vocabulary is limited to the 100,000 objects mentioned by name in the literature, a small fraction of the total known universe. Our graph captures only objects that individual papers mention. The SDSS list millions to billions of sources. The resulting graph therefore reflects objects mentioned in papers, not the full census of known sources.

Future work could extend this foundation along several directions: richer temporal evolution (e.g., graph snapshots), stronger graph-based inference, and what is mentioned by name in the literature.

The concept-object knowledge graph, including the code and study-mode annotations, are available in the Data and Code Availability section.

THE MESSY REALITY OF DATA:

- OCR is imperfect,
- LLMs make mistakes,
- "Discovery" is hard to define.



WHAT'S NEXT? (Because science is never done),

Our graph only includes objects mentioned in papers. That's a tiny fraction of what's out there! Future work could be **MUCH BIGGER.**



DATA AND CODE AVAILABILITY

All code to reproduce the object extraction, entity resolution, training and evaluation pipeline, along with (i) the SIMBAD-resolved object identifier mapping and object metadata used in this work, (ii) paper-level extracted object mentions with rule and study-mode annotations, is available at: <https://github.com/JinchuLi2000/astro-link-forecasting>.

ACKNOWLEDGMENTS

This work was supported in part by the U.S. National Science Foundation under Grant AST-3408729 and by a Humboldt Research Award from the Alexander von Humboldt Foundation (YST). This work also used GPT-5-mini API access provided through the NSF National Artificial Intelligence Research Resource (NAIRR) Pilot.

Work at Argonne National Laboratory was supported by the U.S. Department of Energy, Office of High Energy Physics. Argonne, a U.S. Department of Energy Office of Science Laboratory, is operated by UChicago Argonne LLC under Contract No. DE-AC02-06CH11357. NR is supported by Laboratory Directed Research and Development (LDRD) funding from Argonne National Laboratory.

Want to try it yourself?
GET THE CODE HERE!
(Open science ftw!)



Thanks for reading! Now go find some cool, unexpected connections! (And maybe cite us if you do 😊).

