

Extracting semantic knowledge and observational data for scientific models

(as a service)

Theo Rekatsinas, Shanan Peters, Miron Livny



Semantic knowledge extraction

- **Goal:** create repositories with ***machine-readable data*** about scientific models across different scientific domains
- **Semantic knowledge:** concepts, facts, relationships, entities, and observational data related to scientific phenomena and models
- **Our focus:** identify, extract, and collect semantic knowledge for scientific models from the published literature
- Foundational for model understanding, evaluating, and enhancing scientific models

Equation

$$\delta^{13}\text{C} = (R_{\text{sample}}/R_{\text{standard}} - 1) \times 10^3,$$

Body Text

where R is $^{13}\text{C}/^{12}\text{C}$. The standard is Pee Dee Belemnite limestone that has been assigned a value of $0.0\text{\textperthousand}$. The precisions of $\delta^{13}\text{C}$ determination were less than $0.2\text{\textperthousand}$. POC and PON concentrations were determined using a TCD detector attached to the elemental analyzer.

For Chl a and pheophytin concentrations, POM samples were extracted in the dark for 12 h by 90% acetone, and their concentrations were measured by the fluorometric method (Japan Meteorological Agency, 1970), using a calibrated Turner Designs TD700 fluorometer. In this study, chlorophyll (Chl) was determined as the total pigment including pheophytin. PO₄-P was extracted filtrate by the ascorbic acid–Mo blue method (Strickland and Parsons, 1965), using a Technicon Auto Analyzer.

Equation

Context

Parameter

Values

Section Header

3. Results

Section Header

3.1. Variations in river discharge and riverine POM composition

Body Text

River discharge of the Kiso Rivers changed considerably during the observation period (Fig. 3). Discharge was low ($<500 \text{ m}^3 \text{ s}^{-1}$) until 22 June, and suddenly increased on 24 June (the first flood, $\sim 2000 \text{ m}^3 \text{ s}^{-1}$), reaching a peak flood on 28 June (the second flood, $\sim 3000 \text{ m}^3 \text{ s}^{-1}$). After that, it

during normal discharge. However, the concentration in the Nagara River at high discharge was the same level as that at normal discharge. After discharge, POC concentrations decreased in all rivers. $\delta^{13}\text{C}$ of POM in the Kiso River and the Nagara River varied from $-27.3\text{\textperthousand}$ to $-23.1\text{\textperthousand}$ and from $-29.7\text{\textperthousand}$ to $-25.9\text{\textperthousand}$, respectively. On the other hand, $\delta^{13}\text{C}$ of POM in the Ibi River remained fairly constant (ca. $-30\text{\textperthousand}$). The C/N ratios varied from 7.8 to 22.3 and reached the highest values during high discharge in all rivers.

Table

Table 1

Summary of physical and chemical variables in the Kiso rivers collected at ~ 15 km upstream from the river mouth

	Discharge ($\text{m}^3 \text{ s}^{-1}$)	POC (mg l^{-1})	PON (mg l^{-1})	$\delta^{13}\text{C}$ (\textperthousand)	C/N (mol ratio)
Kiso River					
20 June	155	0.61	0.06	-27.3	12.6
28 June	1257	1.78	0.09	-25.5	22.3
4 July	269	0.30	0.03	-23.1	12.5
Nagara River					
20 June	63	2.28	0.34	-27.7	7.8
28 June	1072	2.11	0.13	-25.9	18.3
4 July	129	0.44	0.06	-29.7	8.7
Ibi River					
20 June	21	1.21	0.14	-30.9	9.8
28 June	622	2.53	0.15	-29.5	20.9
4 July	63	0.60	0.10	-29.0	7.9

Challenge: reasoning about semantic context

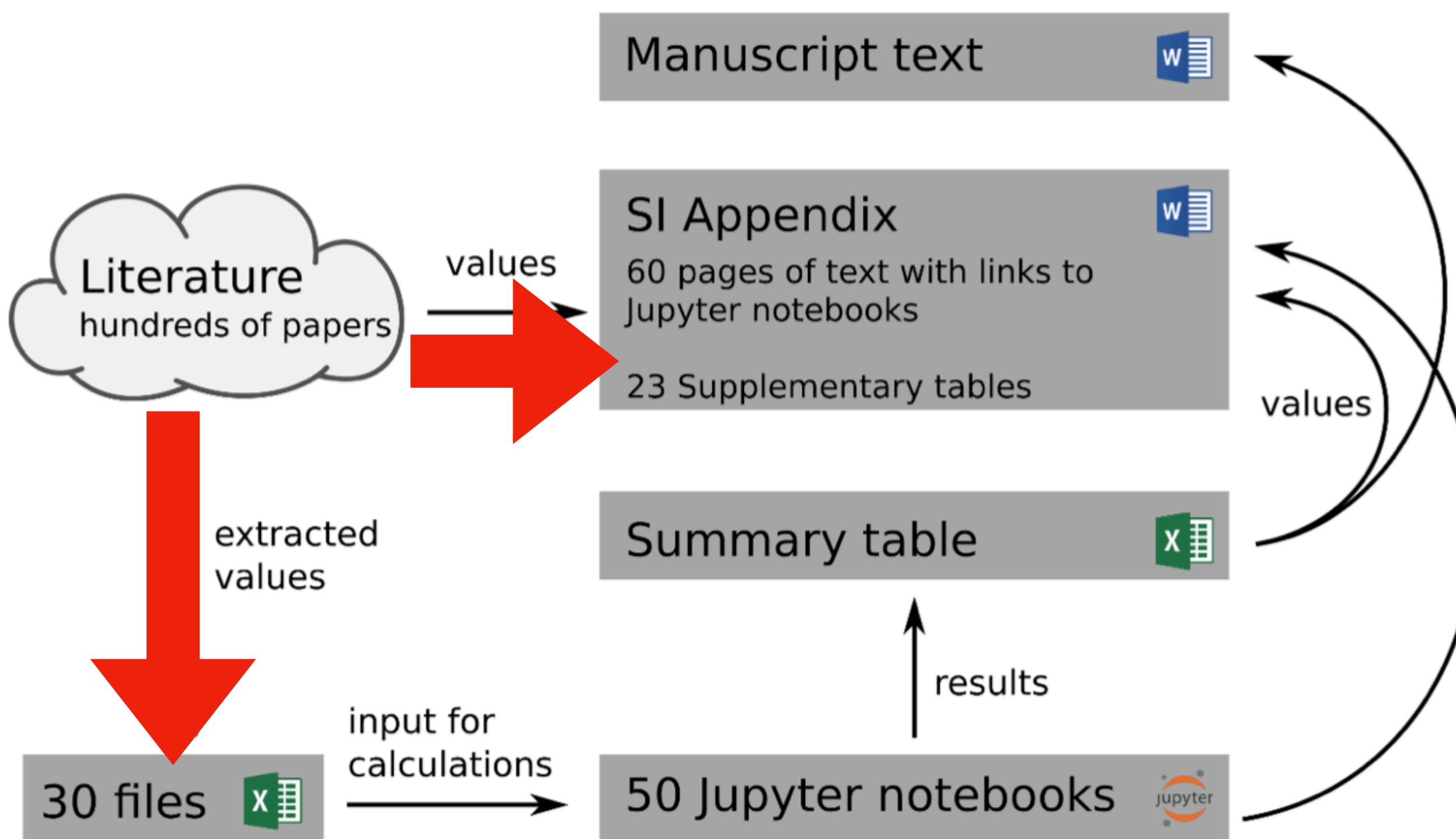
Example: Understanding the biomass distribution on Earth



Proceedings of the
National Academy of Sciences
of the United States of America

The biomass distribution on Earth

Yinon M. Bar-On^a, Rob Phillips^{b,c}, and Ron Milo^{a,1}



Analyzing the biomass on Earth

- The distribution and amount of biomass on Earth is fundamental to many models of carbon cycling and biological response to environmental changes
- Example: 2018 study depended on extracting published measurements of biomass and their context/classification
- Getting more such measurements is critical to building better models of biomass distribution

Automating data collection with xDD and COSMOS

Example use case

- We used xDD to find all document for a list of arthropod taxonomic names; **Result: 8,231 documents**
- Used COSMOS to extract page-level entities and tables; **Result: 30,701 table objects recognized**
- Used COSMOS to extract the content of tables relevant to the context of our analysis; **Result: 10,157 tables extracted**

Table

Independent variables	Insect variables			
	No. chewer species	Specialist ratio	Abundance	Biomass
1) Taxonomic isolation	-0.268	-0.346	-0.183	-0.174
2) Age	0.311	0.317	0.369	-0.198
3) Geographical distribution	0.364	-0.023	0.267	-0.041
4) Altitudinal range	0.697	0.440	0.463	0.531
5) Leaf flush	0.556	0.555	0.283	0.169
6) Young leaves	0.891	0.704	0.678	0.575
7) Leaf expansion	-0.098	-0.355	-0.183	0.230
8) Tree height	0.250	0.121	-0.035	-0.526
9) Leaf mass	0.088	-0.092	-0.251	0.106
10) Abundance	0.524	0.471	0.475	0.745
11) Leaf water	-0.477	-0.290	-0.135	-0.269
12) Leaf nitrogen	-0.279	-0.102	0.112	-0.076
13) Leaf palatability	0.225	-0.003	-0.093	-0.348
14) Log of ant abundance	0.588	0.809	0.645	0.315
15) Sampling effort	-0.211	-0.264	-0.165	0.325

and Basset Yves, *Local Communities of Arboreal Herbivores in Papua New Guinea: Predictors of Insect Variables*, *Ecology*, 77(6), 1996, DOI: 10.2307/2265794

Preview table

Download pickled pandas dataframe

Download OCRed text

See full stored object

COSMOS can accelerate data collection while being robust and domain agnostic

Ongoing exploration: continental groundwater storage

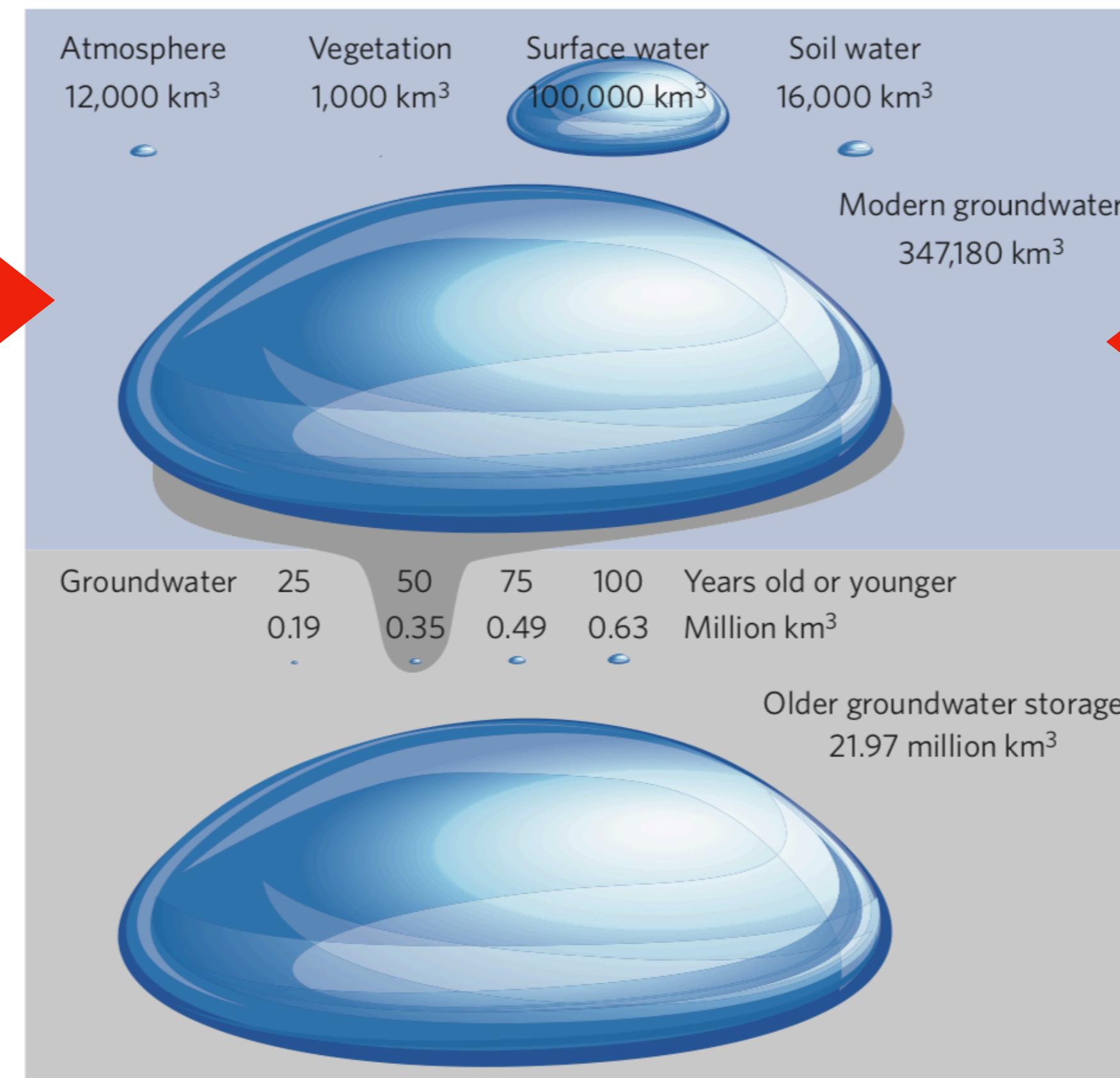
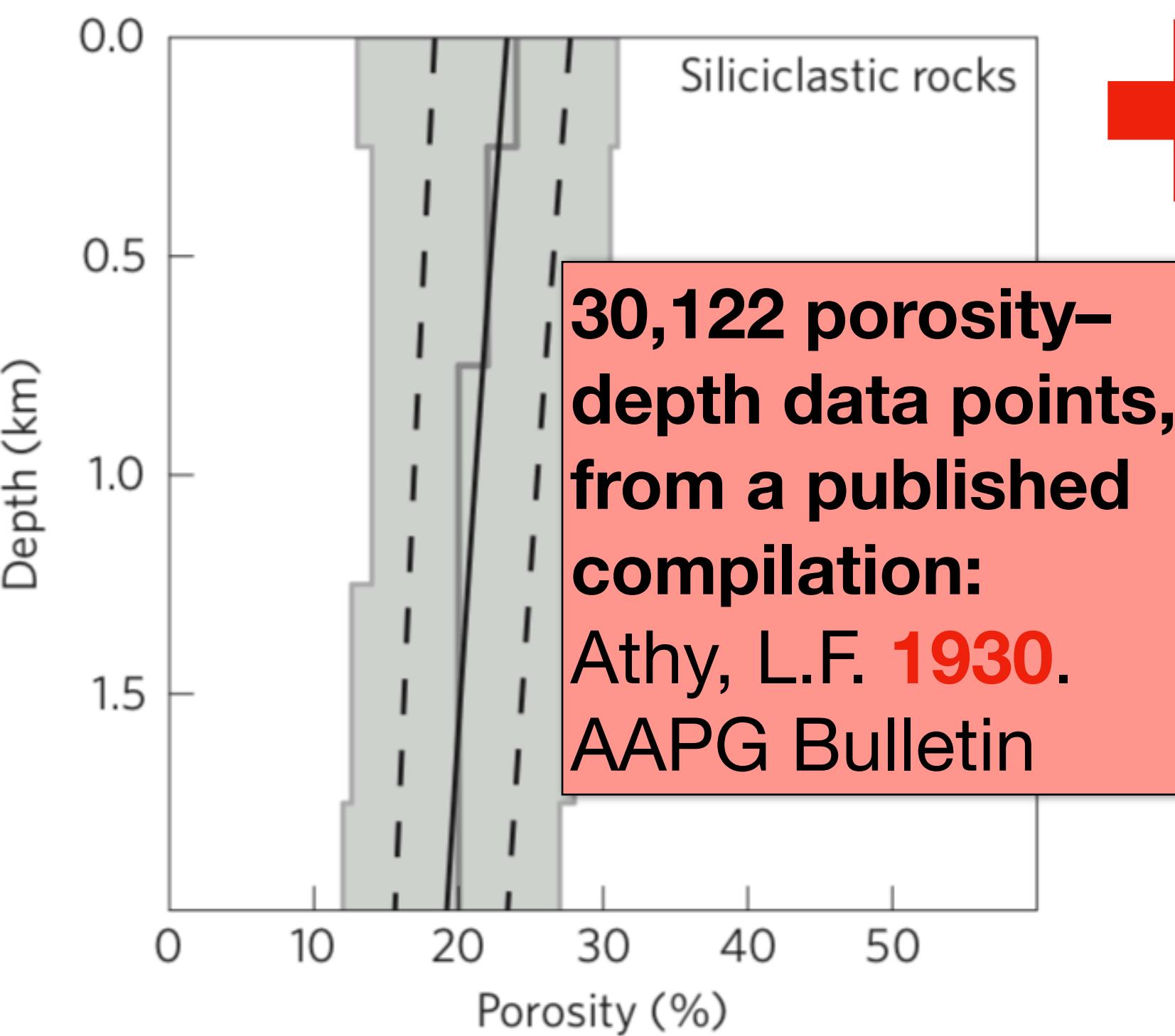
nature
geoscience

Gleeson et al.

ARTICLES

PUBLISHED ONLINE: 16 NOVEMBER 2015 | DOI: 10.1038/NGEO2590

The global volume and distribution of modern groundwater



3,769 measurements of ³H in groundwater compiled from 160 published data sets.

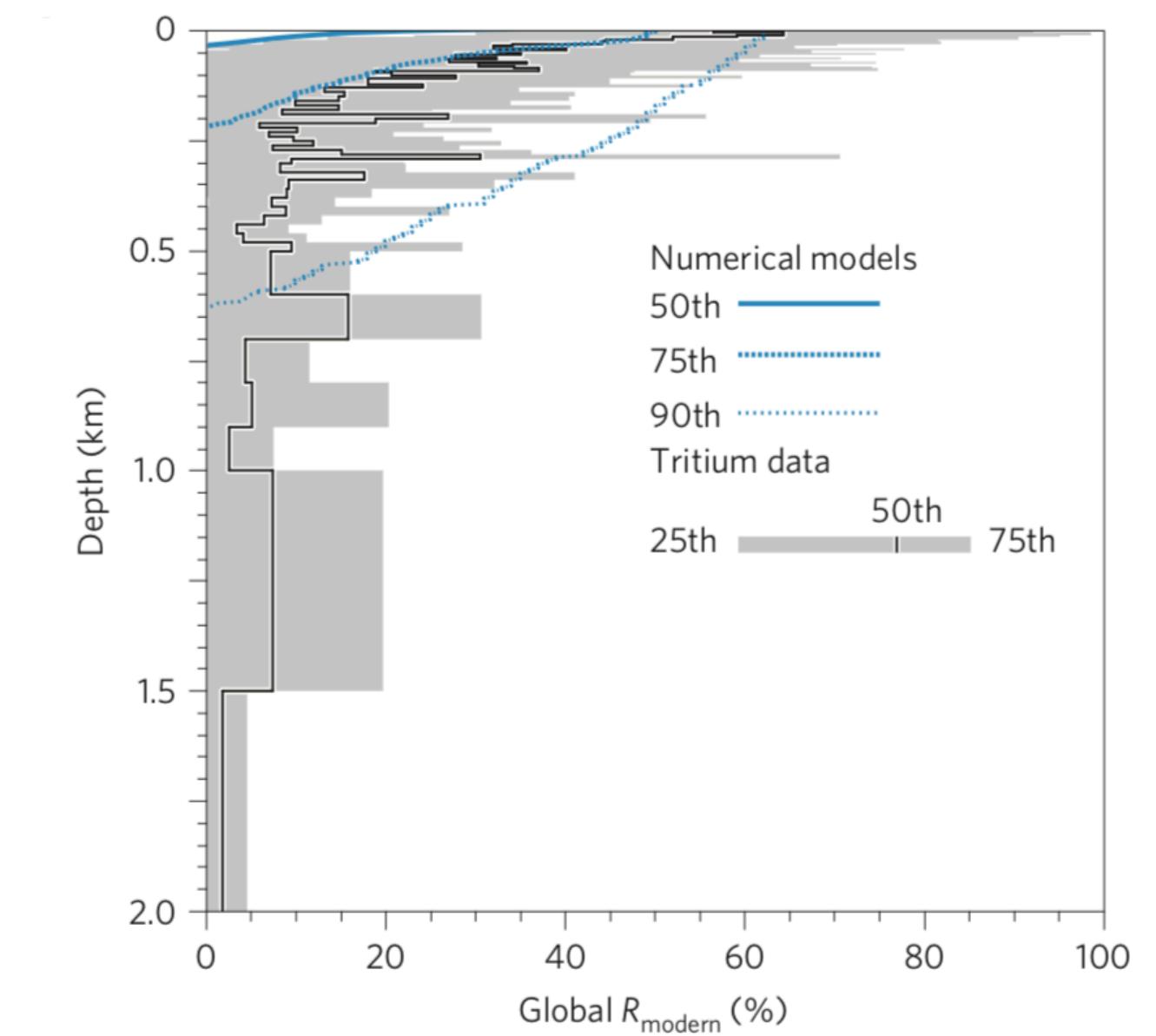
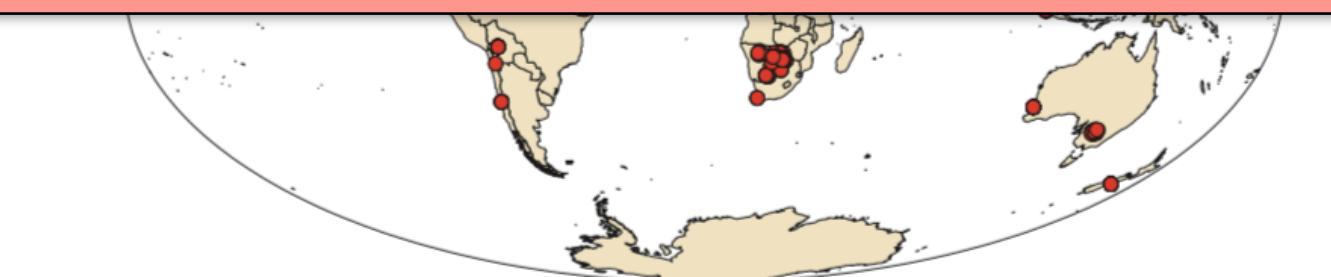


Figure 5 | The different volumes of water stored in the global water cycle.

Ongoing exploration: continental groundwater storage



COSMOS INPUT: NGDS vocabulary for specific measurement types and example context, xDD document acquisition/processing pipeline

Automate location and aggregation of contextualized sample-based measurements and models relevant to groundwater, hydrothermal E



Improve groundwater models and inventories, national hydrothermal assessment

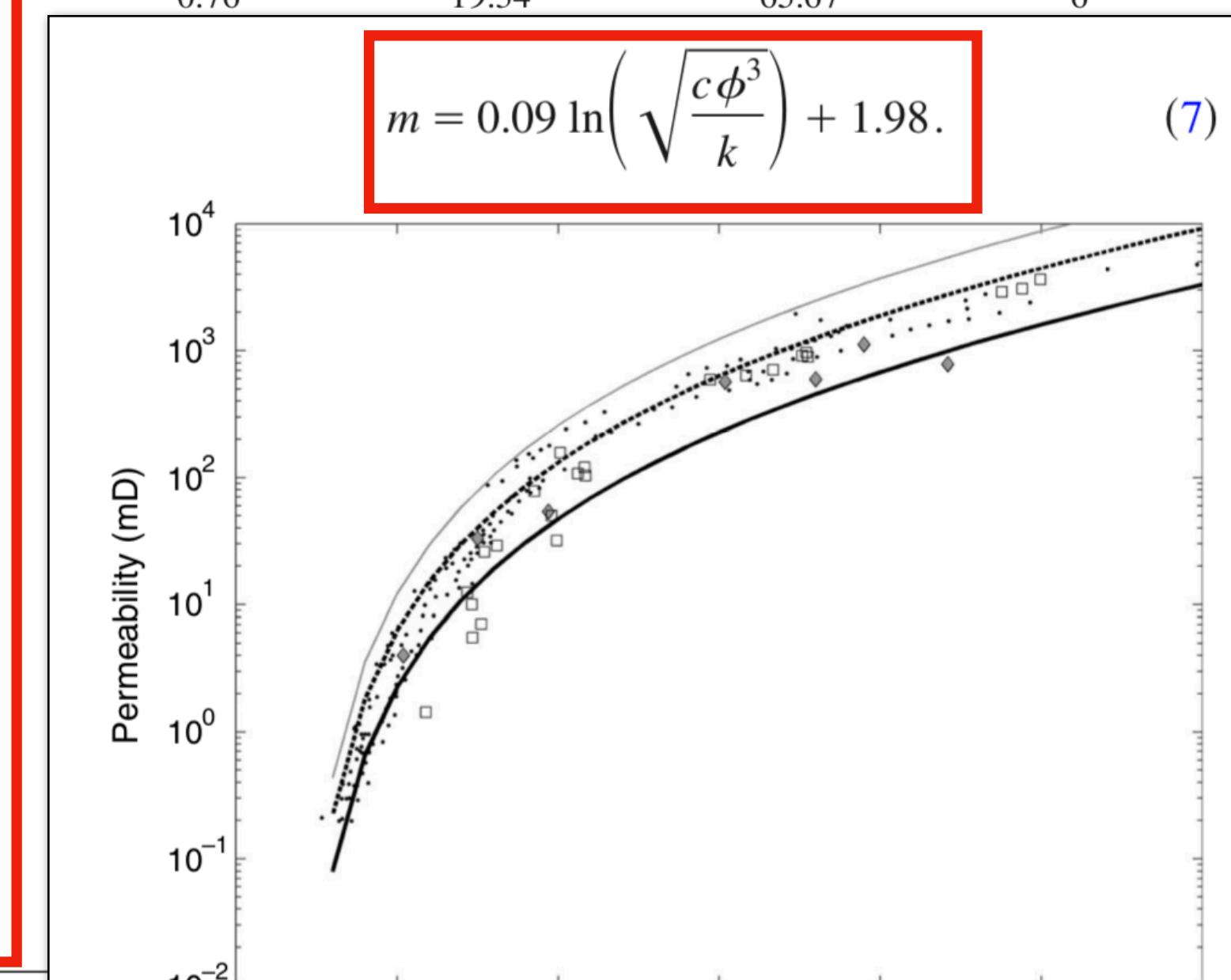
Table 1. Porosity is the helium porosity, permeability R is measured resistivity in ohm·m, F is the normalization factor of 2, and $R_w = 0.17$ ohm·m. Error (%) is the error.

Sample	Porosity	Perm (mD)
A11	0.07	10
A16	0.07	6
A33	0.07	12
A82	0.08	7
A87	0.10	50
A89	0.08	26
A117	0.11	103
B31	0.11	107
B86	0.09	78
B101	0.11	121
B102	0.10	157
B108	0.08	29
F510	0.15	592
GT3	0.17	704
GW18	0.16	637
GW19	0.18	912
GW23	0.18	965
GW28	0.18	896
H27	0.25	3630
H42	0.24	2894
H74	0.24	3079
F410	0.06	1
F570	0.10	32

ABSTRACT

The relations among the resistivity, elastic-wave velocity, porosity, and permeability in Fontainebleau sandstone samples from the Ile de France region, around Paris, France were experimentally revisited. These samples followed a permeability-porosity relation given by Kozeny-Carman's equation. For the resistivity measurements, the samples were partially saturated with brine. Archie's equation was used to estimate resistivity at 100% water saturation, assuming a saturation exponent, $n = 2$. Using

$$m = 0.09 \ln\left(\sqrt{\frac{c\phi^3}{k}}\right) + 1.98. \quad (7)$$



Current status of COSMOS

Retrieve-and-
Read Q&A

On-demand
Knowledge Extraction

Synonym
Services

Dataset
Generation

COSMOS Microservices API

COSMOS
Microservices

Automated container deployment, scaling, and management using Docker swarm; open source



service-oriented

COSMOS
Ingestion Layer

Deep Learning models for document analysis and recognition. Granular information extraction from tables, figures, text, and mathematical expressions.



COSMOS
Knowledge Representation

High Throughput
Computing Layer



GPU, CPU nodes
DAGman job queue
encrypted binaries



parsed and annotated
documents, databases

Digital Library
Layer (xDD)

Automated ingestion of full publications from multiple publishers, backed by institutional agreements: **11.2M documents, +10K daily**



secure storage, backup,
metadata, text indexes

Completed components of COSMOS

- a domain-agnostic ingest engine for unstructured data
- fine-grained entity and data tagging
- services for training natural language comprehension models over user-defined verticals
- multi-modal dataset generation
- analysis and tagging of model code (limited to Fortran)

On-going extensions and research

- On-demand knowledge and data extraction
- Retrieve-and-read Q&A
- Scaling of global-context embeddings
- Powering scientific applications (macroscopic analysis of phenomena)

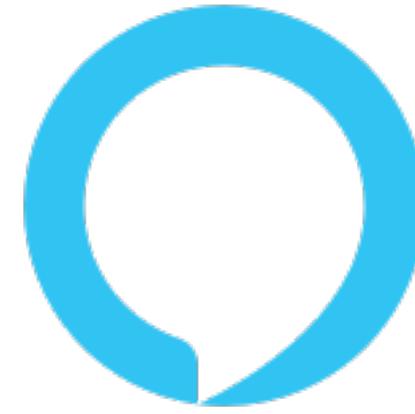


: currently working on it

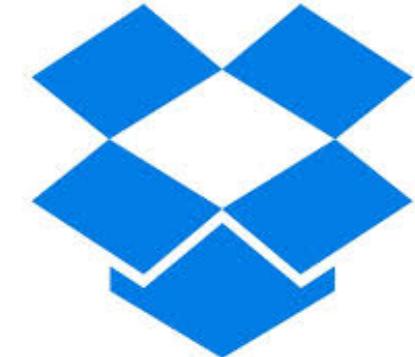


: alpha version completed

Our vision: democratize knowledge extraction



Agent



Applications



Services



Infrastructure

- An intelligent virtual assistant that can consider and reason about semantics across domains when linking/synthesizing artifacts from publications
- APIs that enable users to build custom applications over that assistant
- Services that bring fine-grained knowledge extraction within reach of every scientist—without requiring programming expertise (interactions based on natural language and point-and-click interfaces)
- Compute infrastructure and data to support diverse domains and disciplines

Knowledge extraction solutions are vertical-specific

Current view on Knowledge Extraction:

- **Closed-world extraction:** align to existing entities and attributes, e.g., (ID_Obama, place_of_birth, ID_USA)
- **Closed Information Extraction (IE):** align to existing attributes but extract new entities, e.g., (Theo, place_of_birth, Greece)
- **Open IE:** not limited by existing attributes (or entities)



: production ready



: limited successes for popular verticals

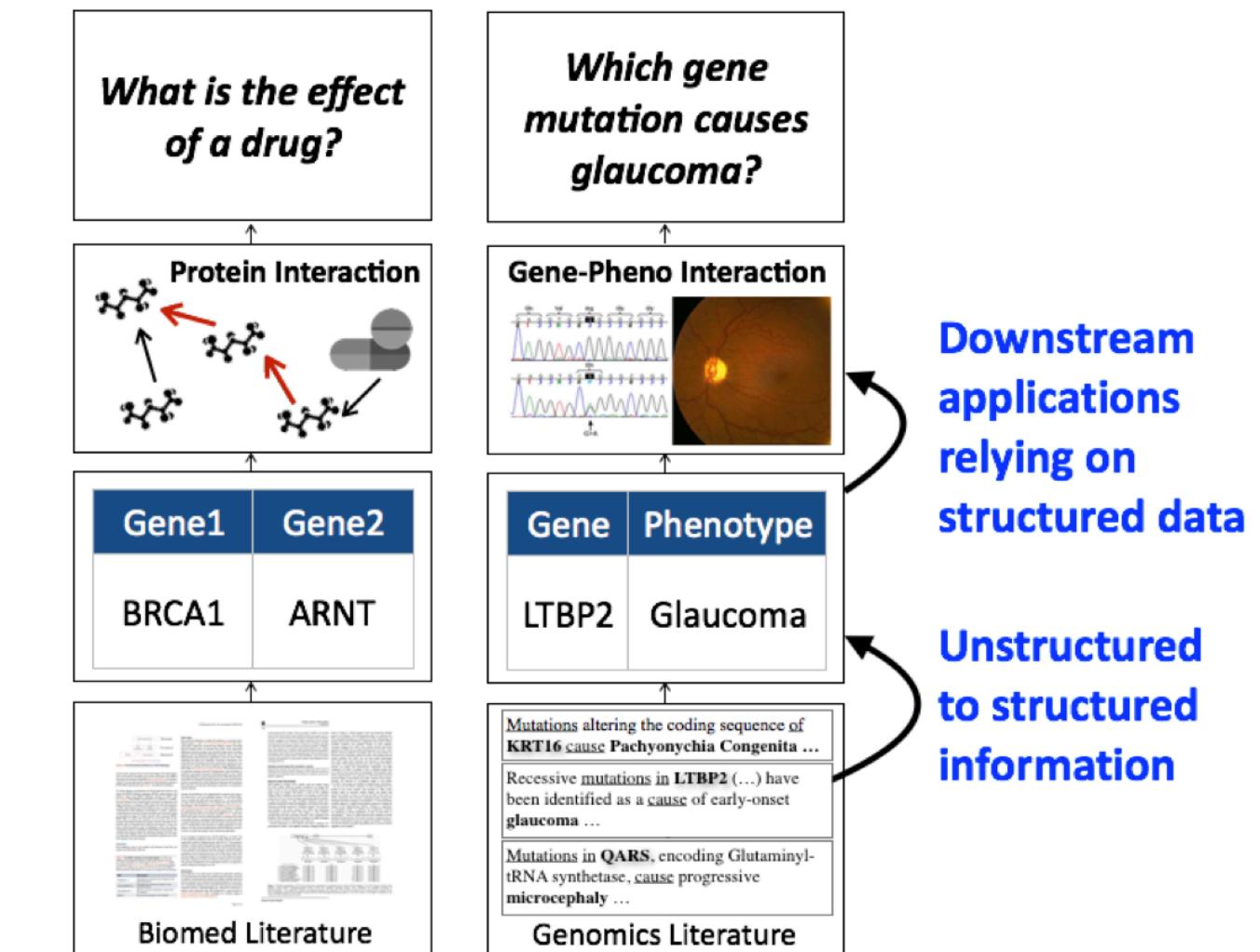
Example systems (from prior DARPA programs):

Memex Human Trafficking



- URL where the advertisement was found
- Phone number of the person in the advertisement
- Name of the person in the advertisement
- Location where the person offers services
- Rates for services offered

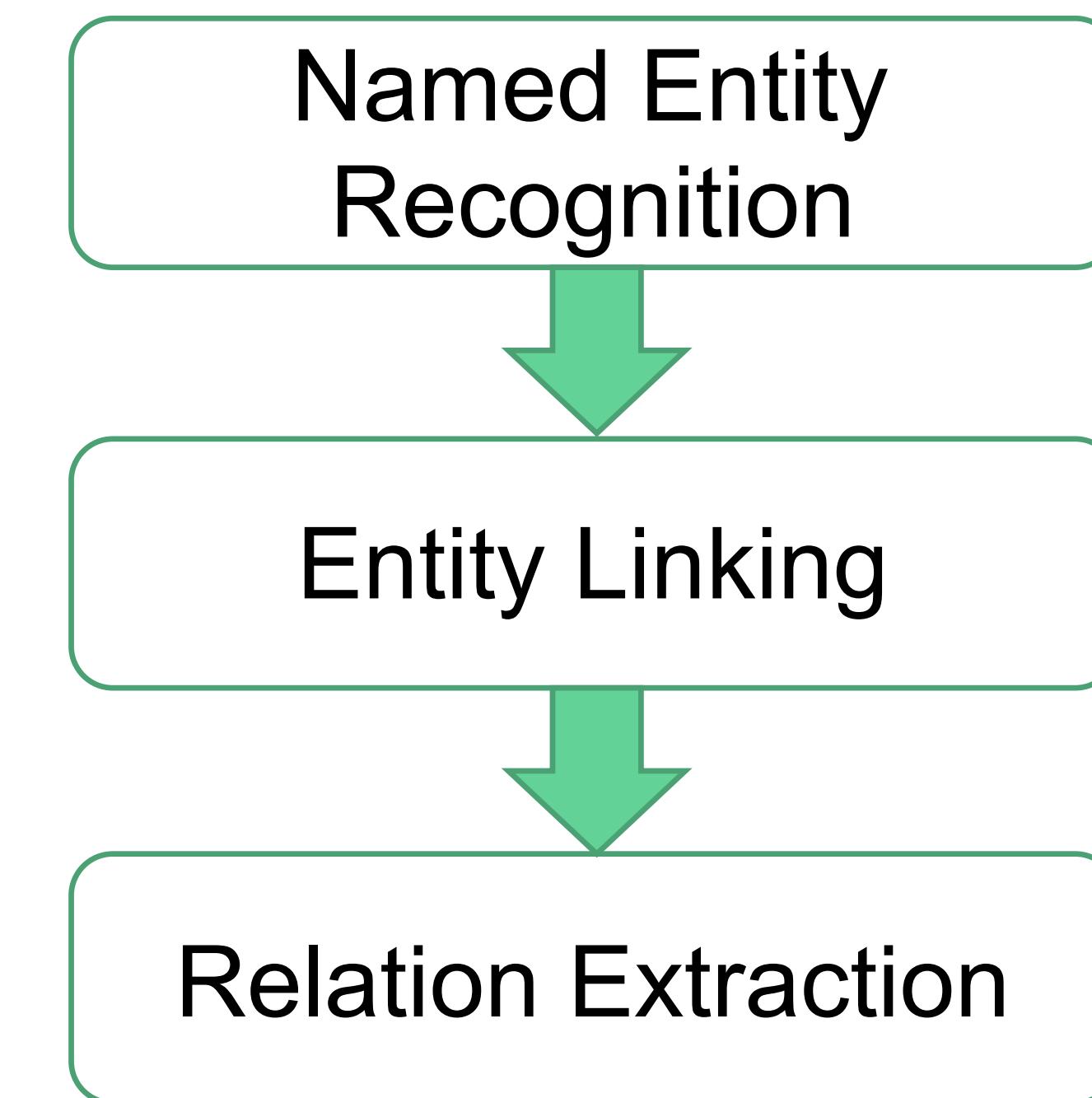
Simplifying Complexity in Scientific Discovery (SIMPLEX)



The bottleneck in Knowledge Extraction is human supervision

State-of-the-art Knowledge Extraction requires supervised machine learning

- **NER:** identify text-spans that are of interest; requires ***matchers or ML models trained*** on corpora of the target vertical
- **EL:** align text-spans to existing canonical entities; requires ***supervised ML models for matching***
- **RE:** align text-spans and identified entities to existing relations; requires ***specialized ML models for a domain schema*** (focus on specific relations in target domain)



Steps in SOTA Knowledge Extraction

The case of unsupervised knowledge extraction with modern text-comprehension models

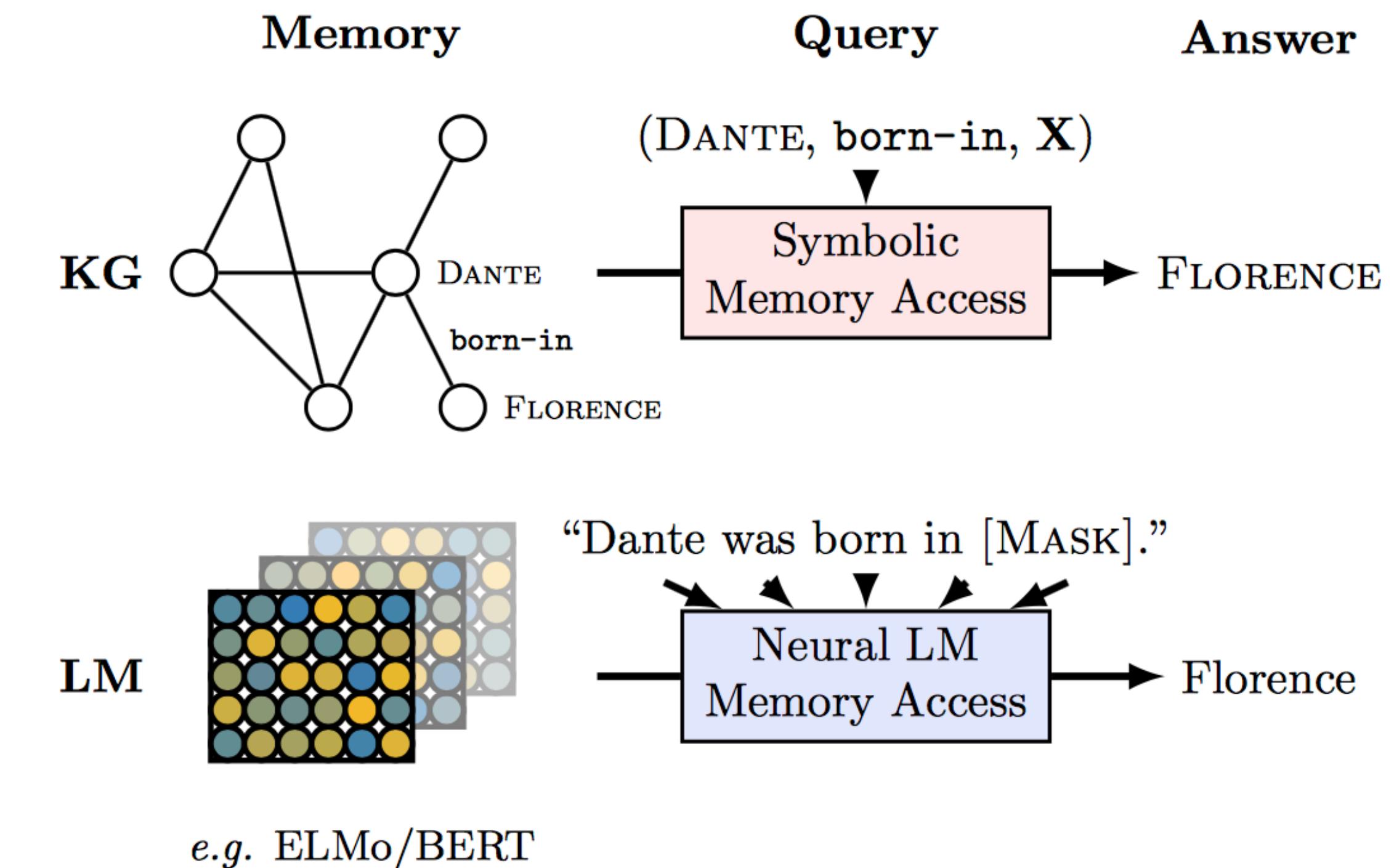
Language Models as Knowledge Bases?

Fabio Petroni¹ Tim Rocktäschel^{1,2} Patrick Lewis^{1,2} Anton Bakhtin¹
Yuxiang Wu^{1,2} Alexander H. Miller¹ Sebastian Riedel^{1,2}

¹Facebook AI Research

²University College London

{fabio petroni, rockt, plewis, yolo, yuxiangwu, ahm, sriedel}@fb.com

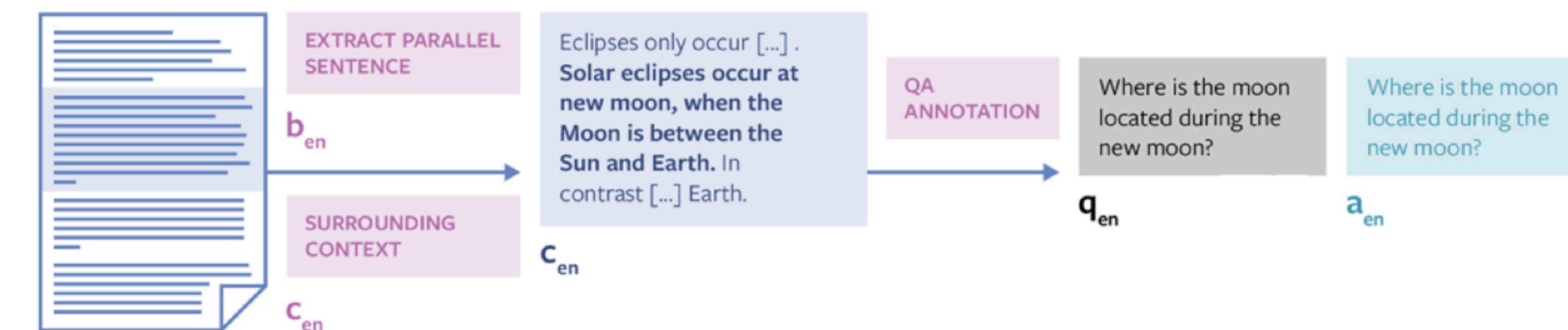
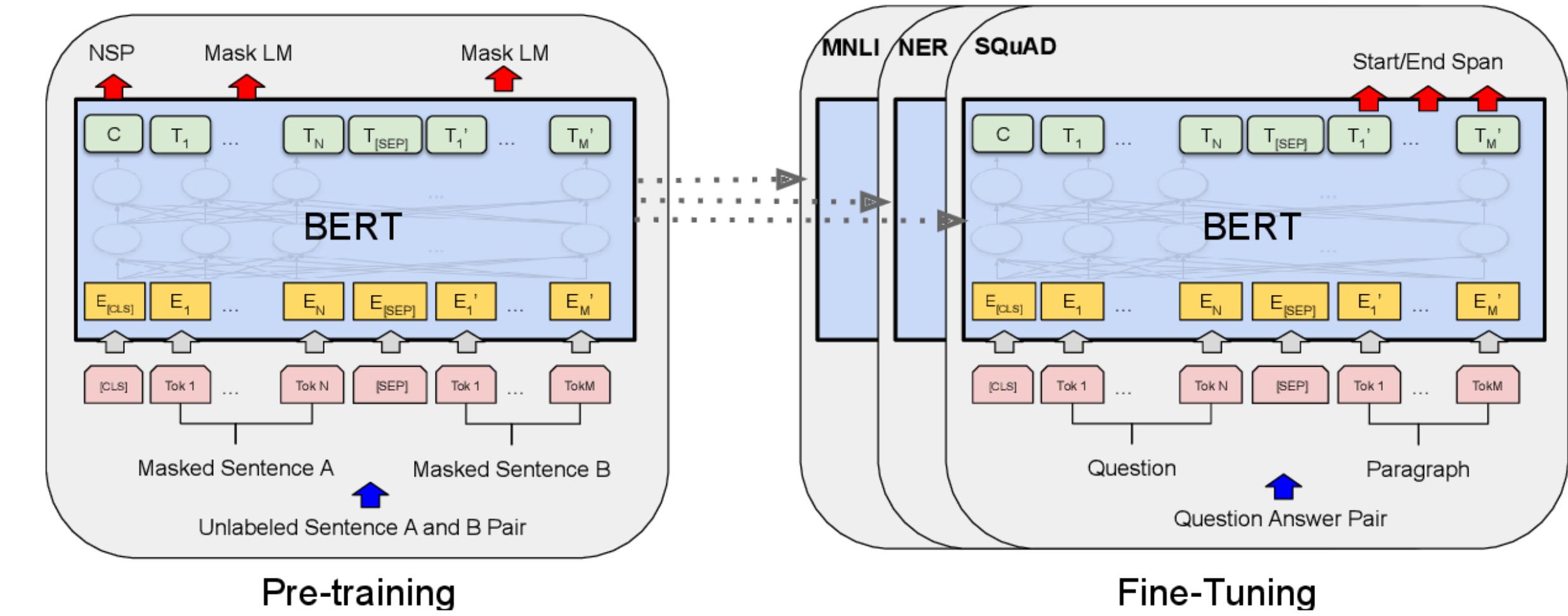


Contextual AI models for text comprehension
(unsupervised models) capture structured knowledge

The case of unsupervised relation extraction with modern text-comprehension models

Limitations:

- Models overfit on heavy-hitters (tend to predict most popular entities); **Cannot be applied to long-tail domains**
- Models need to be fine-tuned to downstream tasks; **Fine-tuning requires training a question-answering head for each domain**
- No end-to-end solution; **QA requires that one has already identified the context that contains a relevant extraction**



Re-Flex: Unsupervised relation extraction with contextual models

Unsupervised relation extraction (Overview)

- **Step 1:** The user provides a dictionary of tokens of interest (subjects); Ex.: set of rock formations, set of people, set of genes
- **Step 2:** The user describes a relation in Natural Language that these tokens participate in; Ex.: located in, plays for, is associated with disease
- **Step 3:** We complete a knowledge base of (subject, relation, object) tuples where the subjects and relations are specified as above. We only use unsupervised, contextual text comprehension models

Input: Components, parts of documents

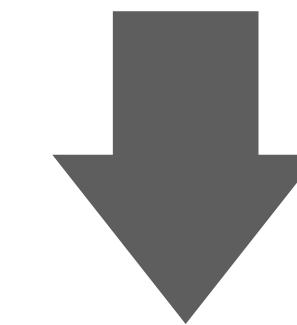
Henrik Lundqvist grew up with his identical twin brother Joel in Åre, Jämtland, an area where ~~skiing~~ is the ~~main~~ sport. Austin Watson was born January 13, 1992, in Ann Arbor, Michigan, where he grew up with his brother, ~~Joel~~, who ~~is~~ a ~~former~~ ice hockey player. ~~He~~ more ~~recently~~ played basketball.

NBA player	Position	Height (inches)	Weight (pounds)
Carmelo Anthony	Forward	80	230
Aaron Baynes	Center	82	260
Matt Bonner	Forward/Center	82	240
Patrick Christopher	Guard	77	209
Norris Cole	Guard	74	170
Glen Davis	Forward/Center	81	289
Gorgui Dieng	Center	83	245
Cleanthony Early	Forward	80	219
Jrue Holiday	Guard	76	205
Jonas Jerebko	Forward	82	231
Andrew Nicholson	Forward	81	250
David West	Forward	81	240

sports.

Input: Text completion queries by combining subjects with relation tokens

Travis Hamonic plays for [MASK]
Morten Madsen plays for [MASK]
Henrik Lundqvist plays for [MASK]



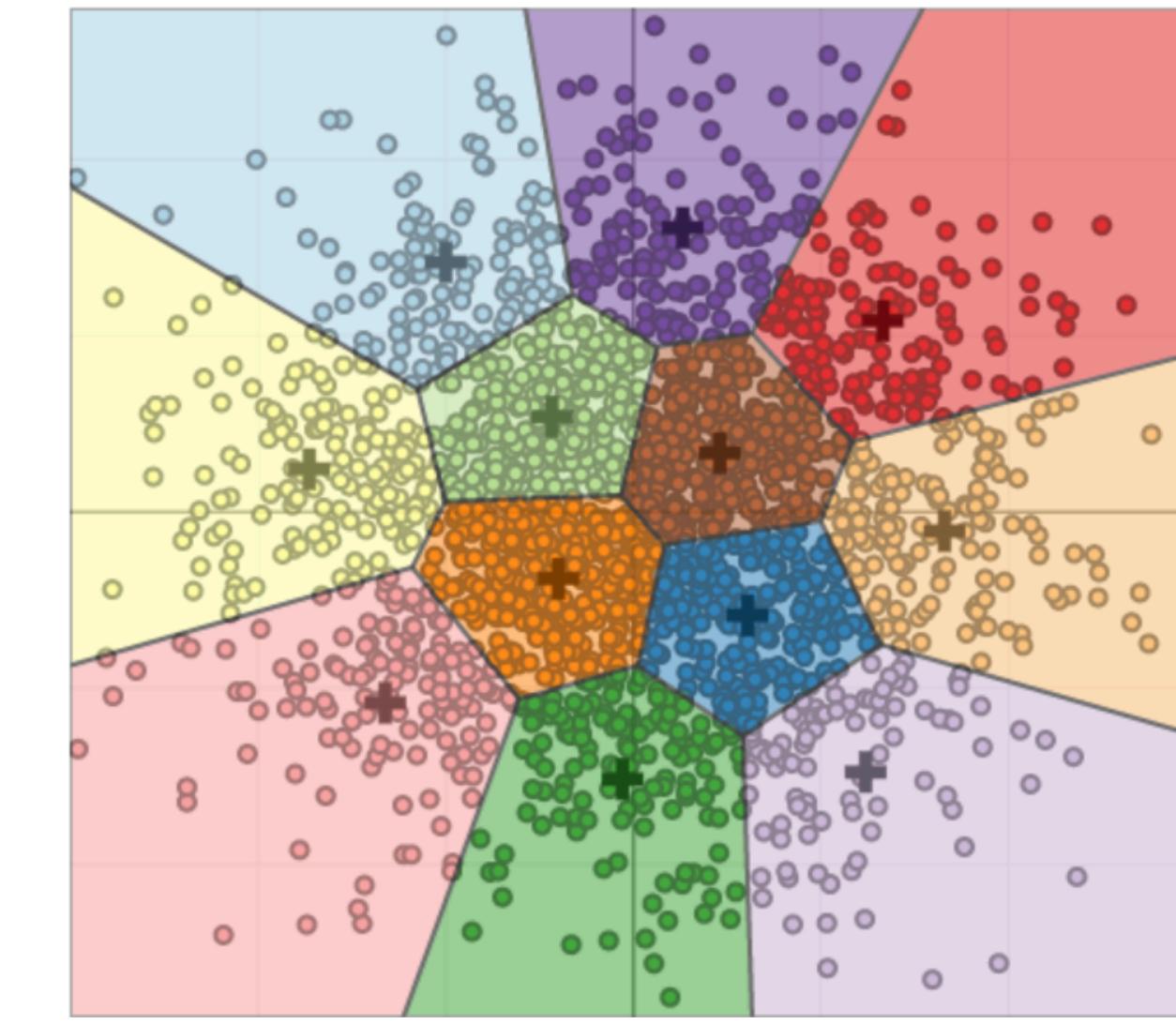
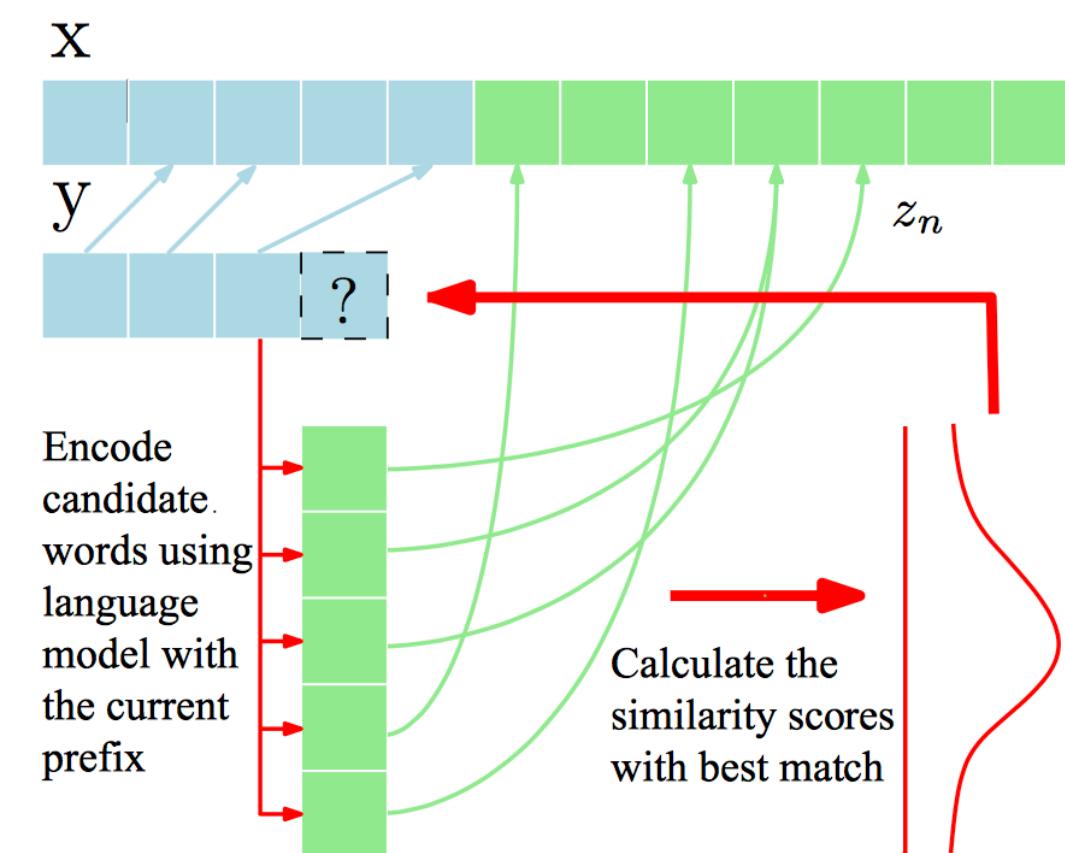
Re-Flex
Step 1: Align queries to relevant contexts
Step 2: Find most probable anchor answer in context (inference fine-tuning)
Step 3: Answer decoding

Output: Completed knowledge base

Travis Hamonic plays for **Calgary flames**
Morten Madsen plays for **Modo Hockey**
Henrik Lundqvist plays for **New York Rangers**

Inference fine-tuning by contextual matching

Condition on context tokens and re-weight their probabilities by computing their contextual similarity to probability distribution over the raw prediction of the model



We use Voronoi clustering to re-assign the probability mass (center correspond to tokens in the relevant context)

We perform inference-level fine-tuning by matching context to raw model predictions

Re-Flex: Unsupervised relation extraction with contextual models

RE-Flex achieves comparable or better performance than **supervised models**

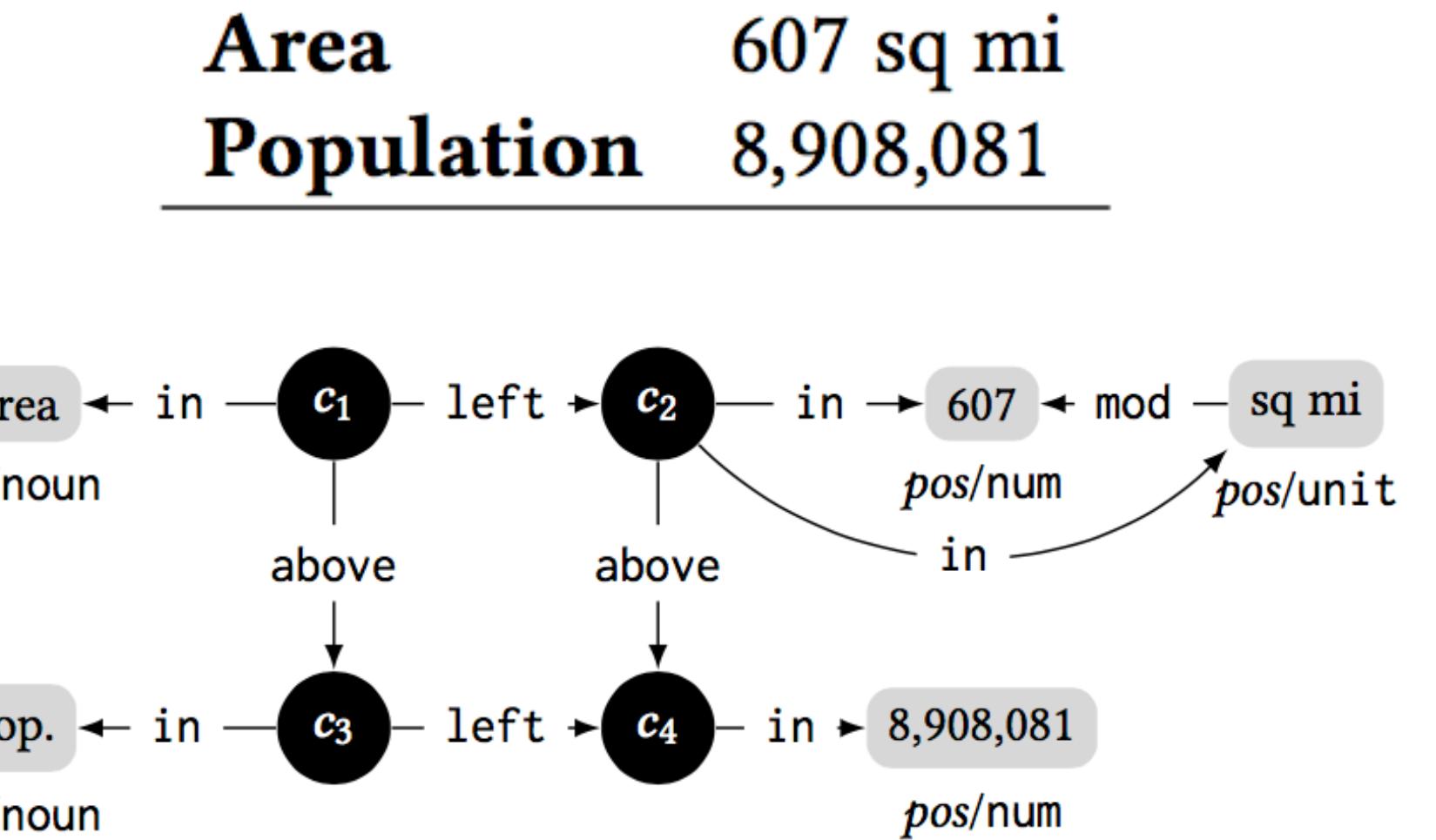
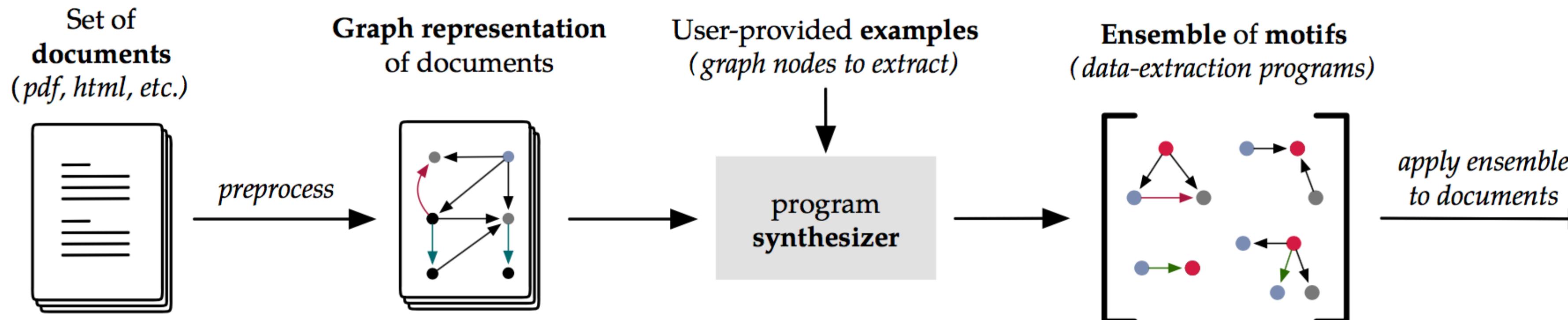
RE-Flex yields ~40 F1 points improvement against SOTA unsupervised methods

Dataset	Model	EM	F1
Google-RE	RE-Flex (Lewis et al., 2019)	91.1	91.1
	(Dhingra et al., 2018)	54.3	68.2
	BERT-Large+Squad2.0	25.9	34.6
	BERT-Large+ZSRE	53.6	77.0
	BERT-Large+ZSRE	83.2	88.2
T-REx	RE-Flex (Lewis et al., 2019)	67.1	67.1
	(Dhingra et al., 2018)	13.6	18.5
	BERT-Large+Squad2.0	20.8	27.7
	BERT-Large+ZSRE	41.2	51.0
	BERT-Large+ZSRE	46.8	52.2
ZSRE	RE-Flex (Lewis et al., 2019)	37.8	41.4
	(Dhingra et al., 2018)	9.6	14.1
	BERT-Large+Squad2.0	12.3	17.4
	BERT-Large+ZSRE	39.4	48.2
	BERT-Large+ZSRE	54.7	56.7

Ensemble Synthesis for Robust Data Extraction

Interactive extraction over COSMOS generated datasets

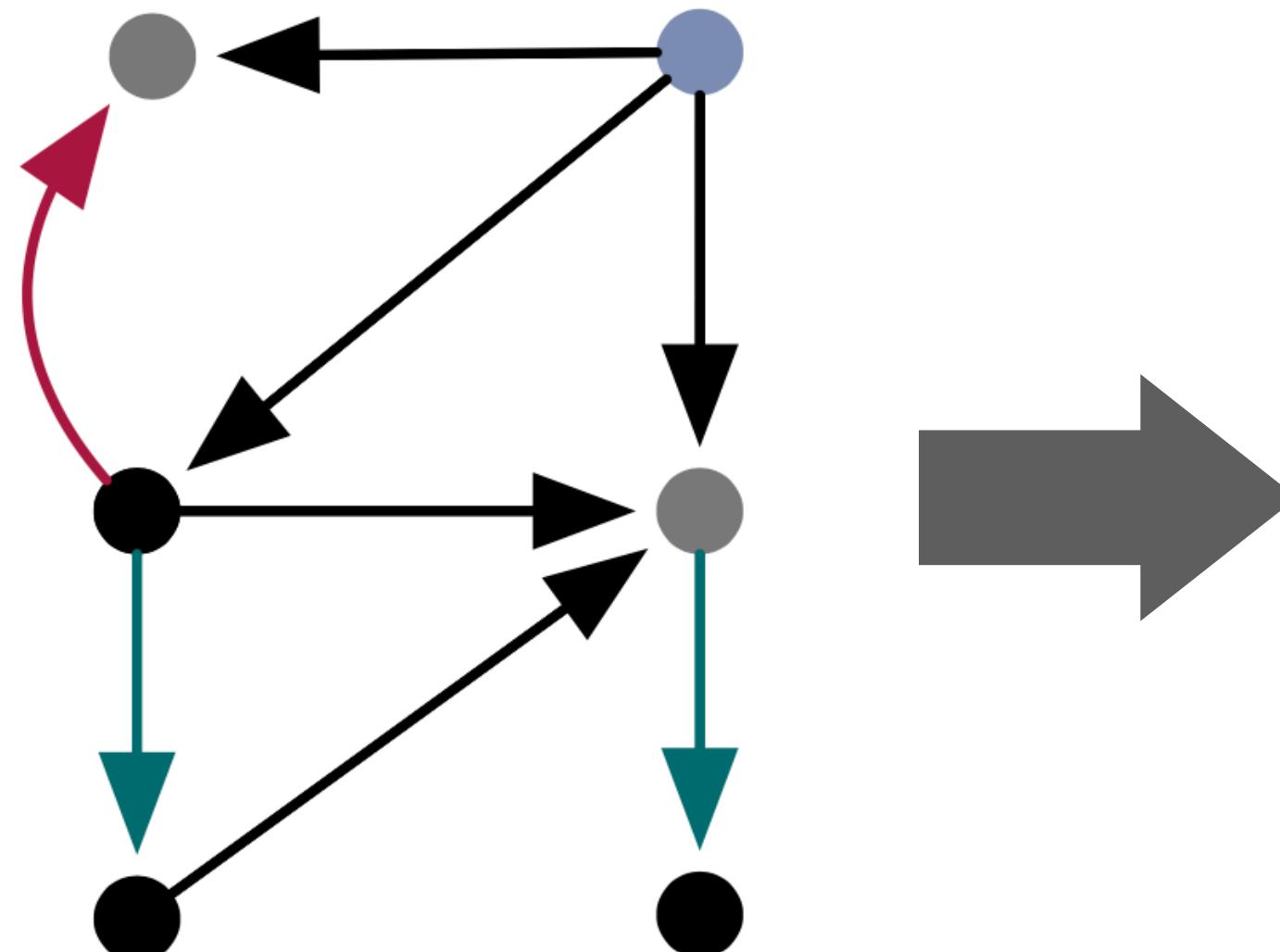
- **Step 1:** Collect related document elements (paragraphs, tables, figures)
- **Step 2:** The user annotates few spans of interest
- **Step 3:** Exploit COSMOS's graph representation and apply program synthesis to generate graph-based programs that extract relevant data from the collection of Step 1
- **Step 4:** Iterate in an active learning loop to improve quality



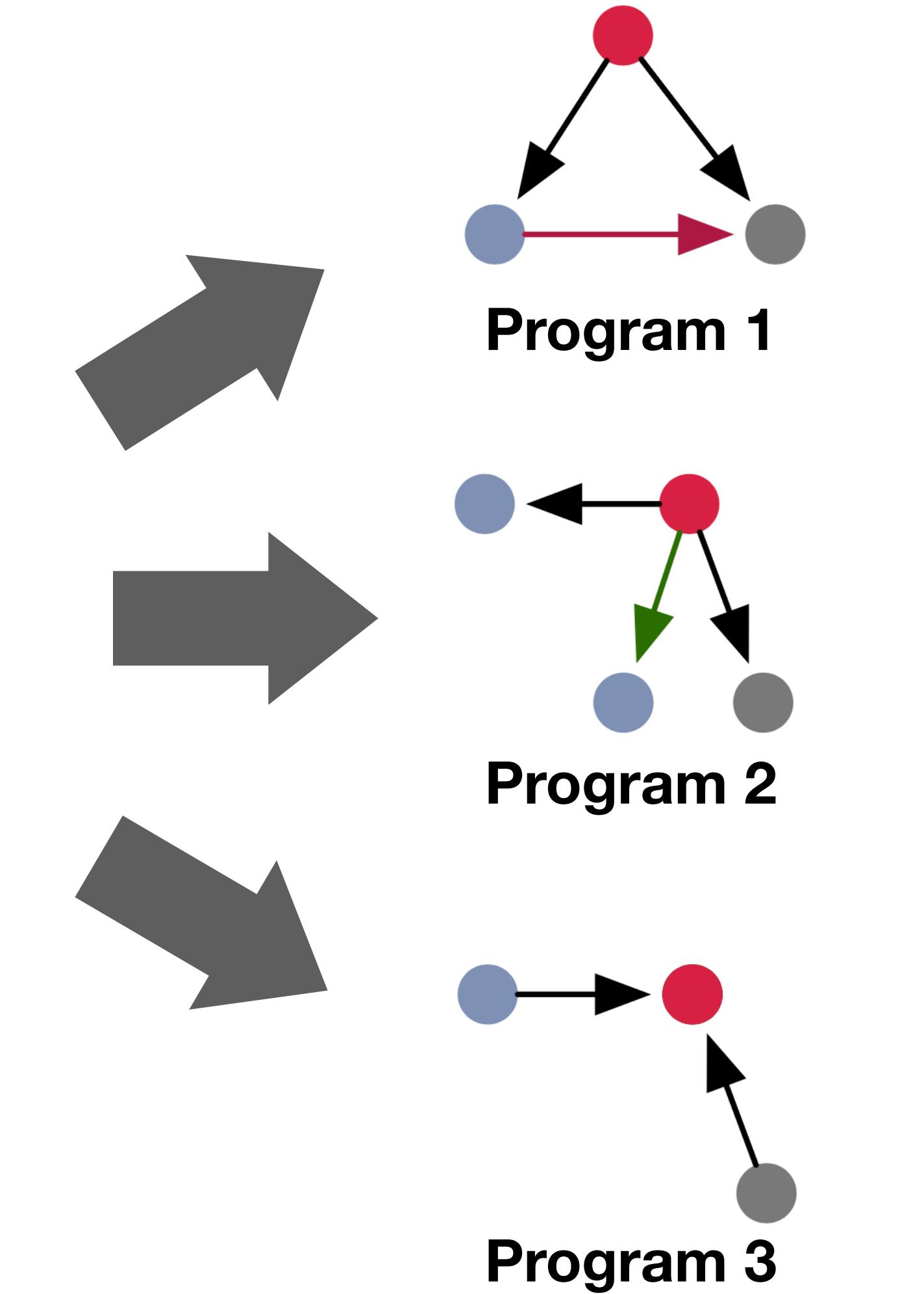
Example Graph Representation of a Table

Program synthesis with Bayesian Ensembles

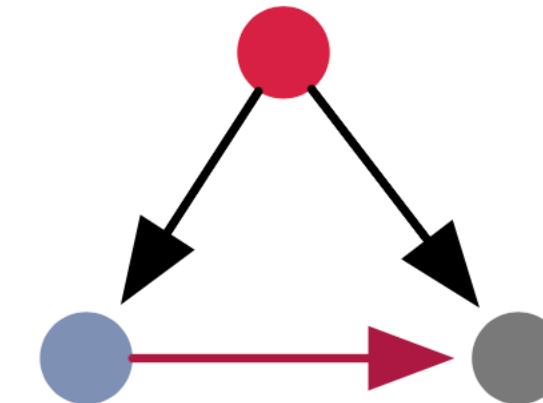
Input: Document graphs
(*unlabeled data*)
limited annotations by
expert (less than 10)



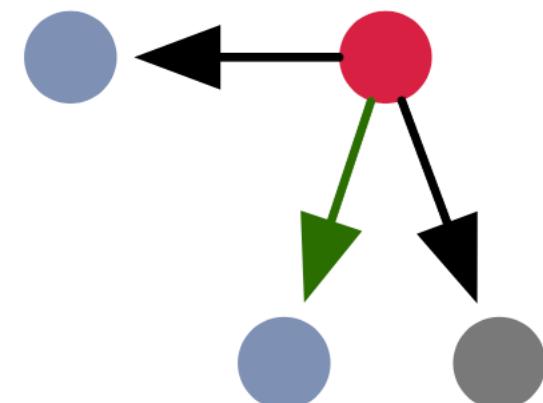
Program Synthesizer
A new program synthesis approach
that generate multiple programs that
satisfy the input specification



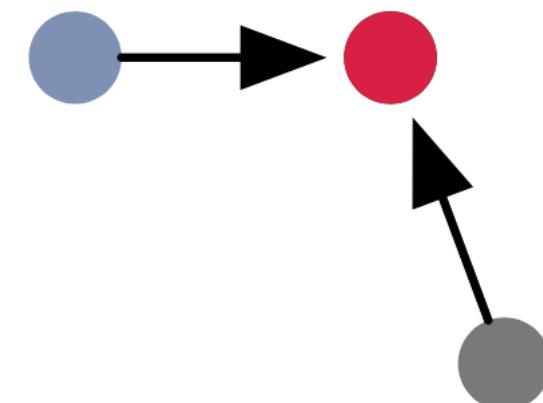
Program synthesis with Bayesian Ensembles



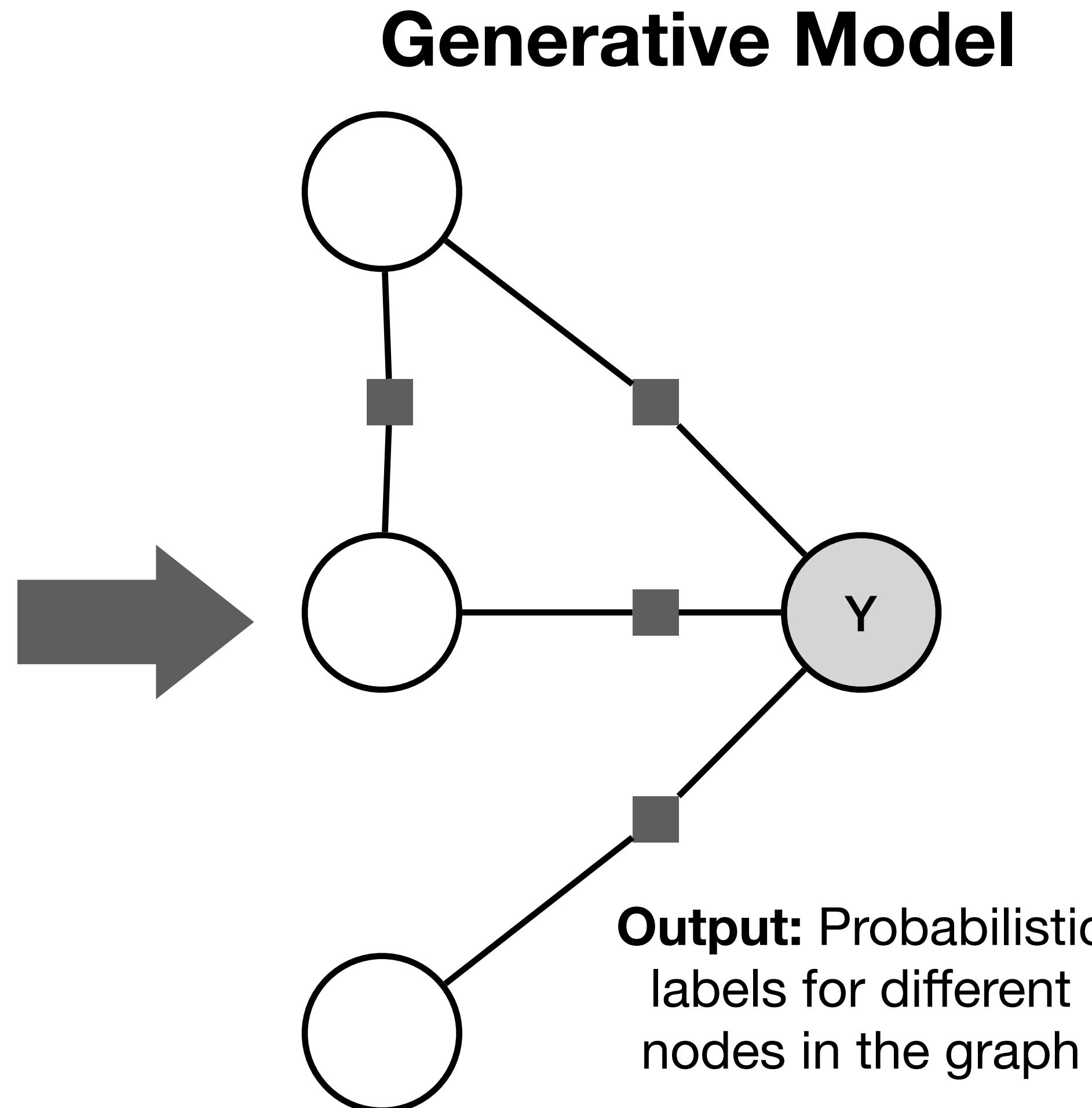
Program 1



Program 2



Program 3



$$Pr[y = 1 \mid \mathcal{M}(v)] \propto \exp \left(\sum_{M \in \mathcal{M}} \ln(\mathcal{A}(M)) \cdot \mathbb{1}[M(v) = 1] \right)$$

Active learning loop

initialize $FPR_0(M)$ for each motif $M \in \mathcal{M}$

initialize learning rate η

for i in $0, \dots, l$

update w_M for each motif M using $FPR_i(M)$

query user for label y of v with the **highest entropy** (see text)

if $y = 1$, then **for** $M \in \mathcal{M}$ s.t. $M(v) = 1$ (true positives)

$FPR_{i+1} \leftarrow FPR_i(M) \cdot \exp(-\eta)$

if $y = 0$, then **for** $M \in \mathcal{M}$ s.t. $M(v) = 1$ (false positives)

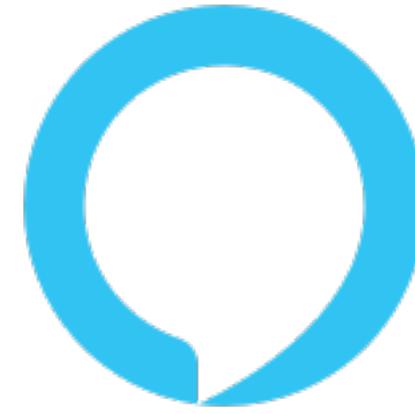
$FPR_{i+1} \leftarrow FPR_i(M) \cdot \exp(\eta)$

Active learning can help improve the accuracy estimate of each program

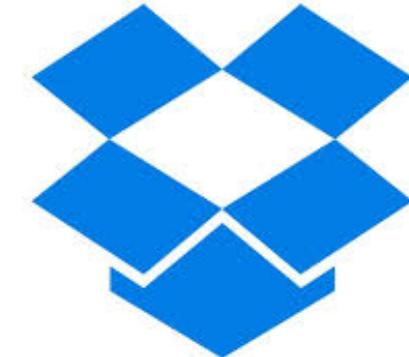
Dataset name	F1-score	# Ex.
Political (Cell)	1.0	1
Chemical (Span)	0.703	5
President (Cell)	0.66	1
Transistor (Cell)	0.65	1
Transistor (Span)	0.5	3

Competing or SOTA performance with less than five examples

Our vision: democratize knowledge extraction



Agent



Applications



Services



Infrastructure

- An intelligent virtual assistant that can consider and reason about semantics across domains when linking/synthesizing artifacts from publications
- APIs that enable users to build custom applications over that assistant
- Services that bring fine-grained knowledge extraction within reach of every scientist—without requiring programming expertise (interactions based on natural language and point-and-click interfaces)
- Compute infrastructure and data to support diverse domains and disciplines

Supplemental Slides

Current status of COSMOS

Retrieve-and-
Read Q&A

On-demand
Knowledge Extraction

Synonym
Services

Dataset
Generation

COSMOS Microservices API

COSMOS
Microservices

Automated container deployment, scaling, and management using Docker swarm; open source



service-oriented

COSMOS
Ingestion Layer

Deep Learning models for document analysis and recognition. Granular information extraction from tables, figures, text, and mathematical expressions.



COSMOS
Knowledge Representation

High Throughput
Computing Layer



GPU, CPU nodes
DAGman job queue
encrypted binaries



parsed and annotated
documents, databases

Digital Library
Layer (xDD)

Automated ingestion of full publications from multiple publishers, backed by institutional agreements: **11.2M documents, +10K daily**



secure storage, backup,
metadata, text indexes

Completed components of COSMOS

- a domain-agnostic ingest engine for unstructured data
- fine-grained entity and data tagging
- services for training natural language comprehension models over user-defined verticals
- multi-modal dataset generation
- analysis and tagging of model code (limited to Fortran)

On-going extensions and research

- On-demand knowledge and data extraction
- Retrieve-and-read Q&A
- Scaling of global-context embeddings
- Powering scientific applications (macroscopic analysis of phenomena)



: currently working on it

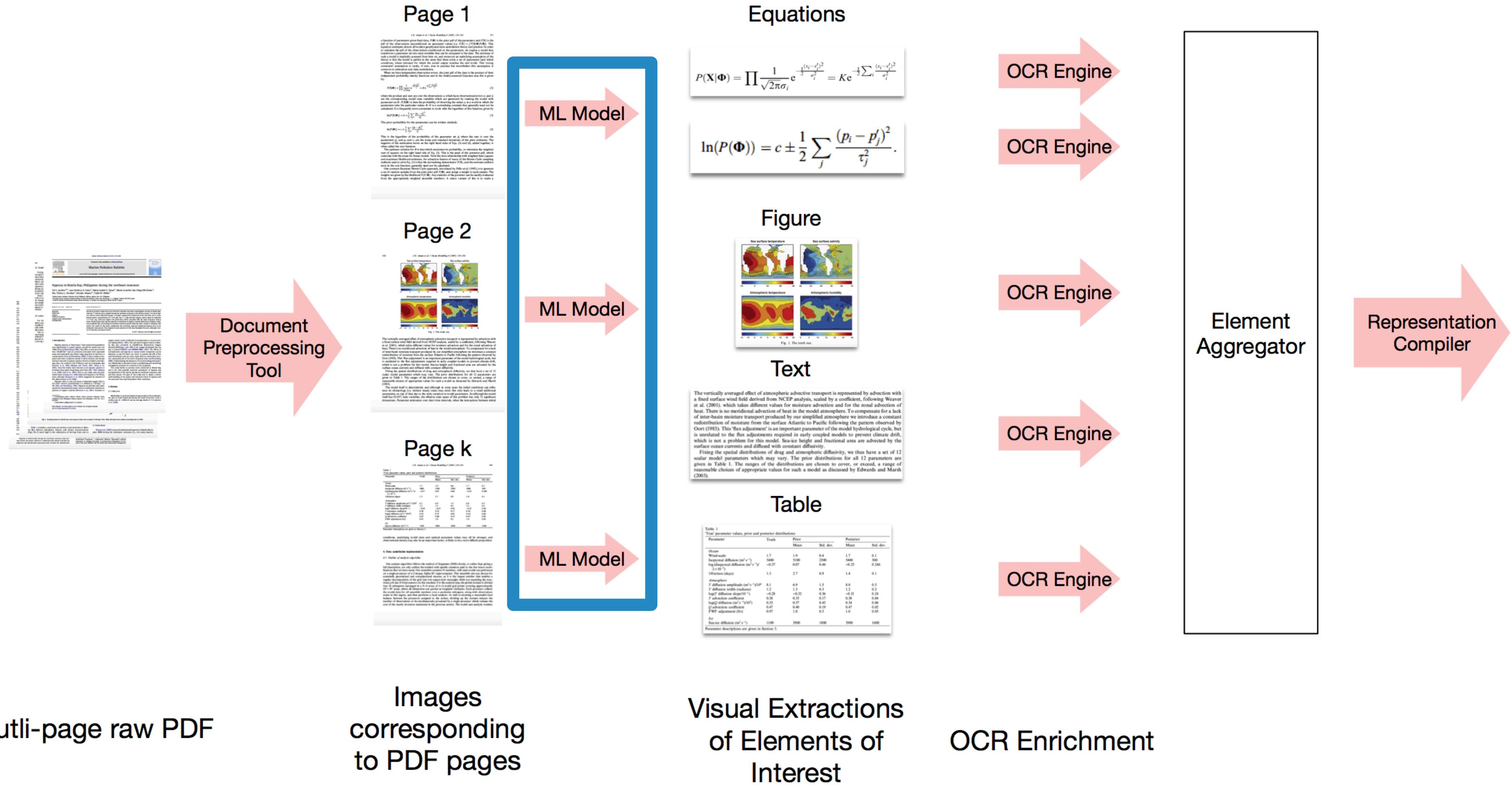


: alpha version completed

From PDFs to XML representations

COSMOS

Ingestion Layer



HTML Representation of initial PDF following Fonduer's unified data model

Images
corresponding
to PDF pages

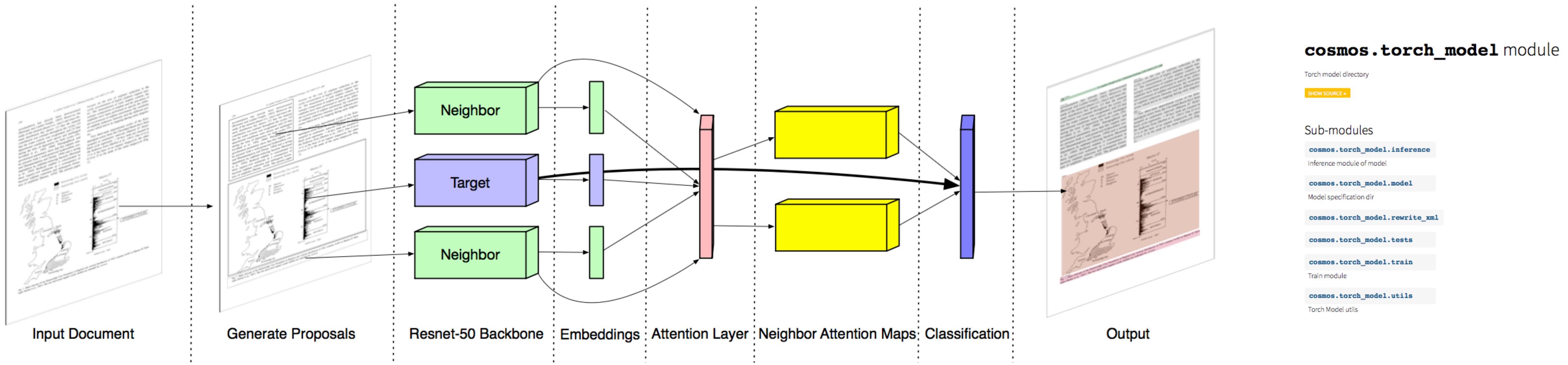
Visual Extractions of Elements of Interest

OCR Enrichment

COSMOS' document segmentation model offers robustness against format heterogeneity across publications

The COSMOS Attentive RCNN Model

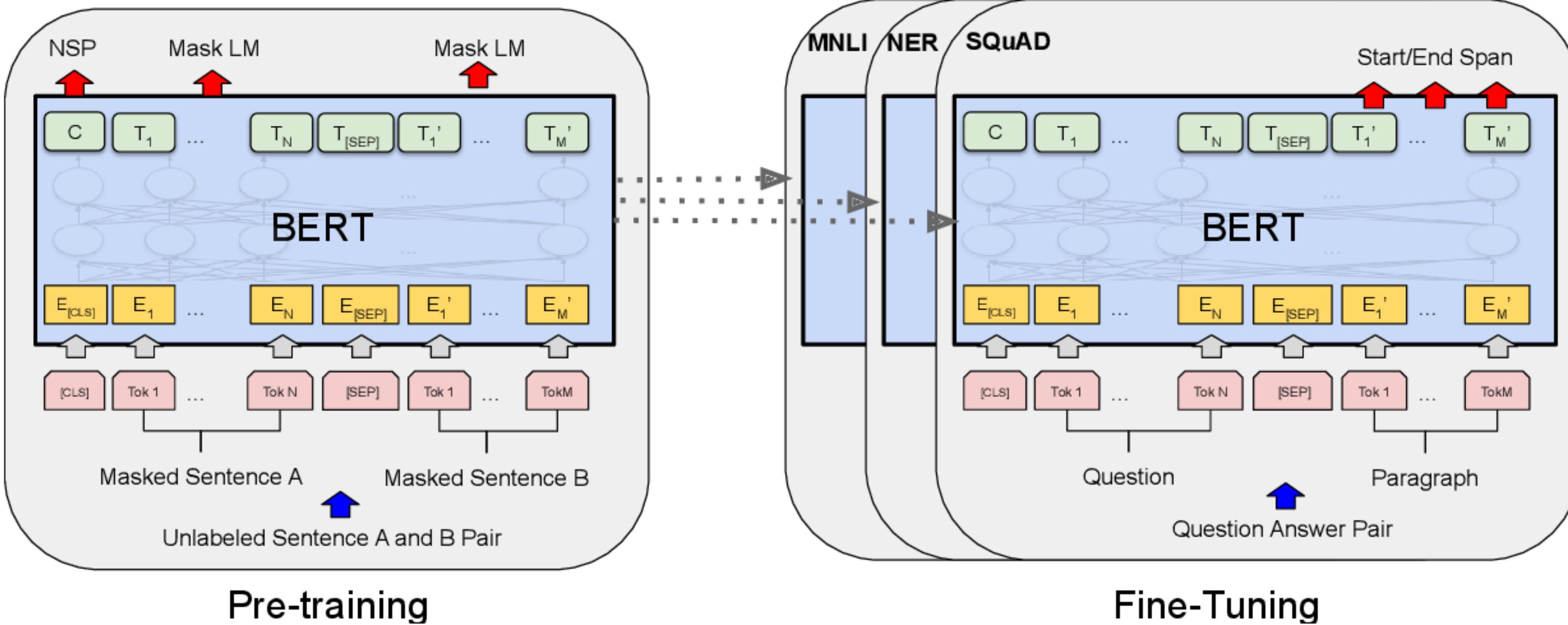
COSMOS
Ingestion Layer



New distributed representation (in the visual space)
for each element in the page.



Data representations in COSMOS



Leaderboard

SQuAD2.0 tests the ability of a system to not only answer reading comprehension questions, but also abstain when presented with a question that cannot be answered based on the provided paragraph. How will your system compare to humans on this task?

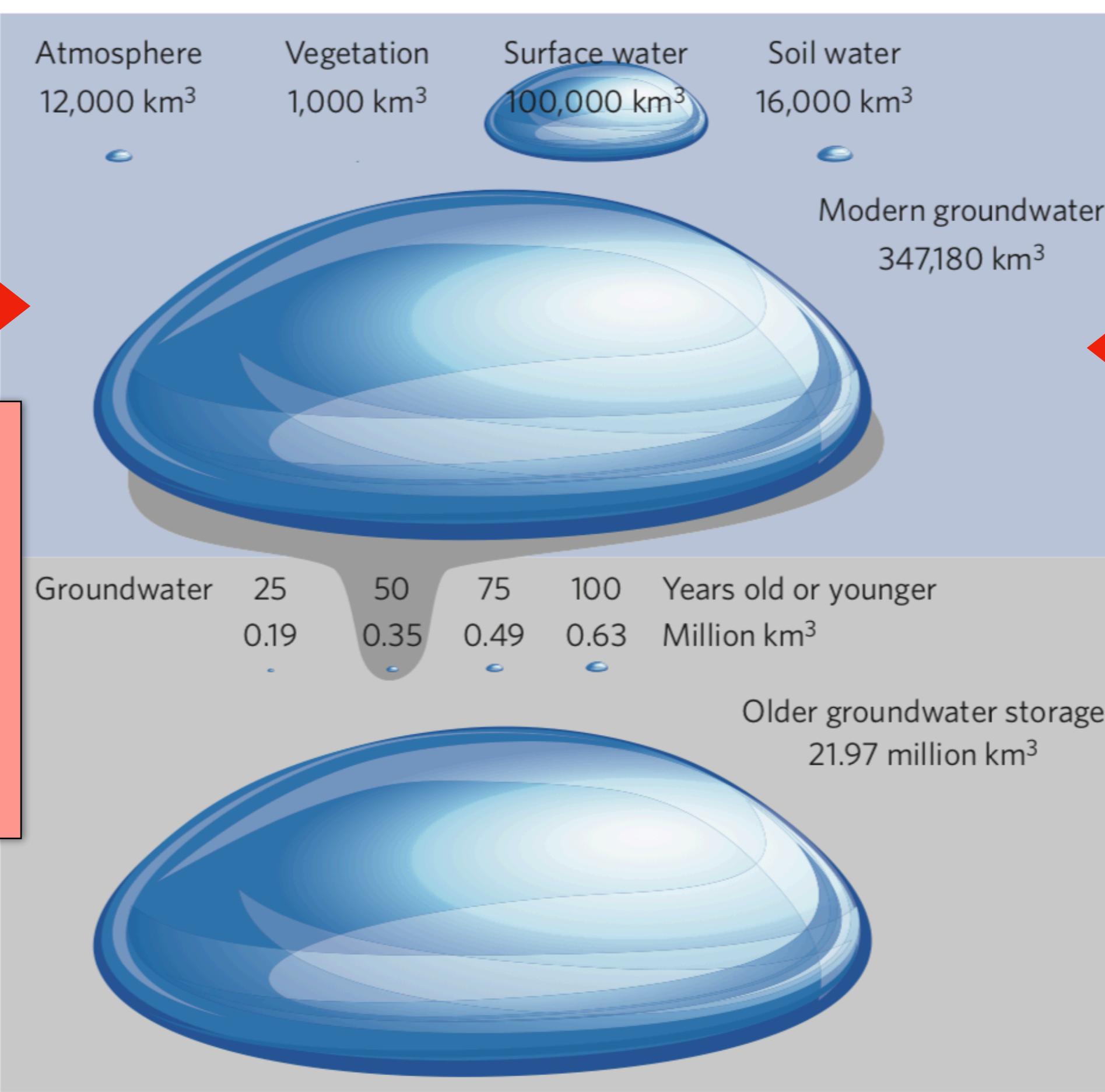
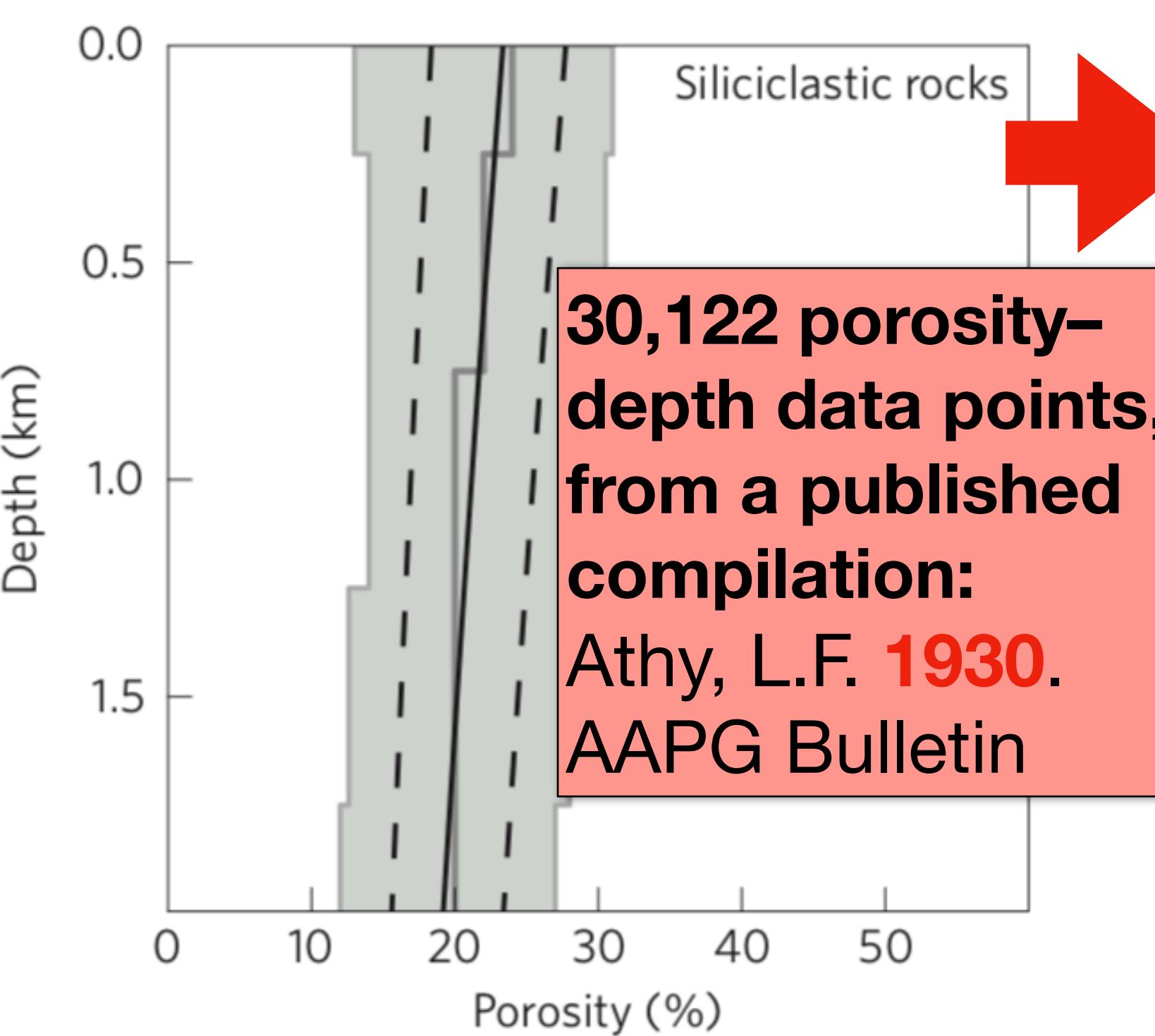
Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1	BERT + DAE + AoA (ensemble) Joint Laboratory of HIT and iFLYTEK Research Mar 20, 2019	87.147	89.474
2	BERT + ConvLSTM + MTL + Verifier (ensemble) Layer 6 AI Mar 15, 2019	86.730	89.286
3	BERT + N-Gram Masking + Synthetic Self-Training (ensemble) Google AI Language https://github.com/google-research/bert Mar 05, 2019	86.673	89.147
4	BERT + DAE + AoA (single model) Joint Laboratory of HIT and iFLYTEK Research Mar 16, 2019	85.884	88.621
5	BERT + MMFT + ADA (ensemble) Microsoft Research Asia Jan 15, 2019	85.082	87.615

We adopt existing entity tagging based representations (e.g., fine-grained named entity tagging and indexing) and state-of-art context-aware representations based on deep learning models (transformer networks)

Use cases that COSMOS can contribute to

Use Case 1: continental groundwater storage

The global volume and distribution of modern groundwater



3,769 measurements of ³H in groundwater compiled from 160 published data sets.

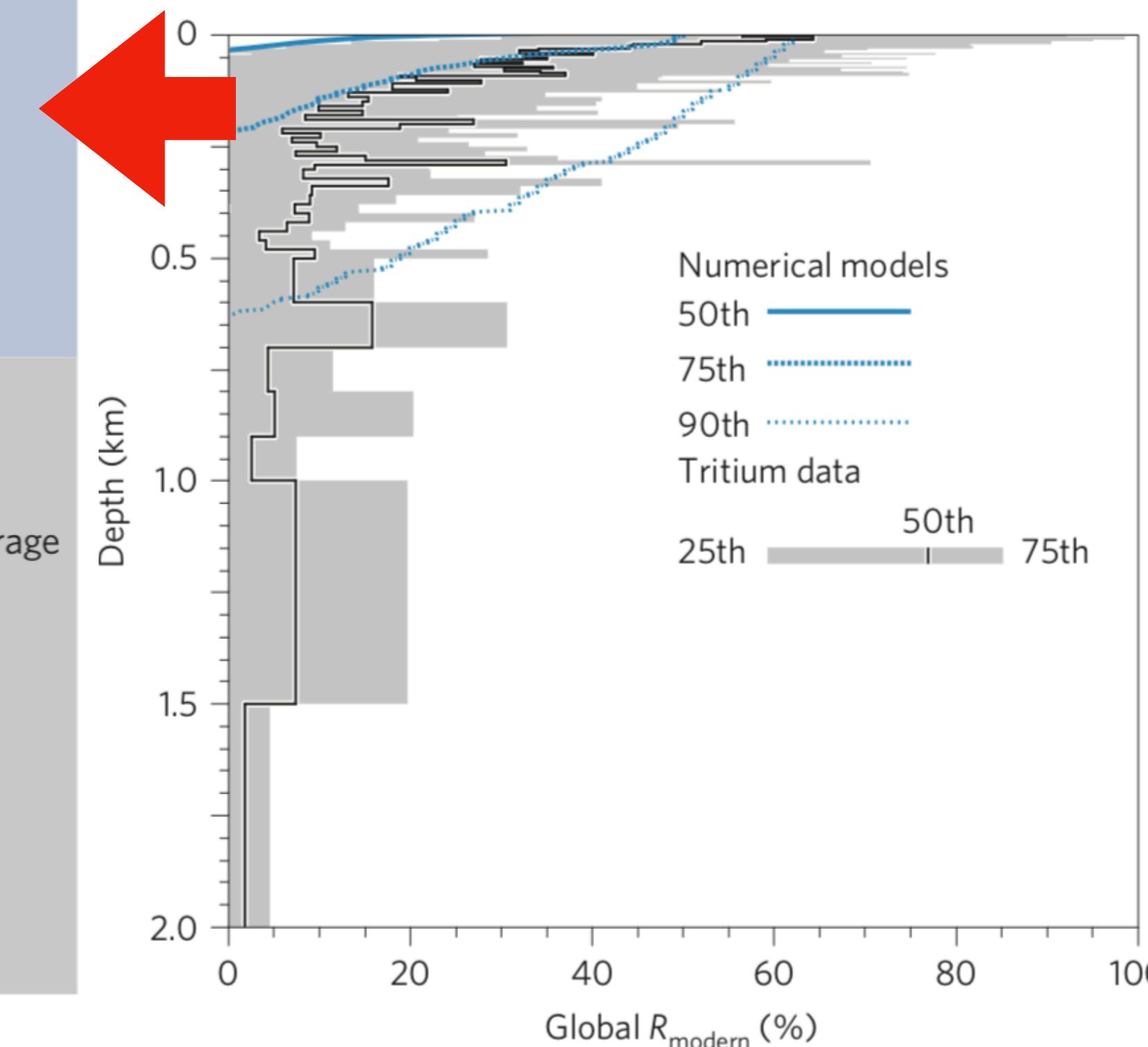


Figure 5 | The different volumes of water stored in the global water cycle.

Use Case 1: continental groundwater storage



Automate location and aggregation of contextualized sample-based measurements and models relevant to groundwater, hydrothermal E

Improve groundwater models and inventories, national hydrothermal assessment

COSMOS INPUT: NGDS vocabulary for specific measurement types and example context, xDD document acquisition/processing pipeline

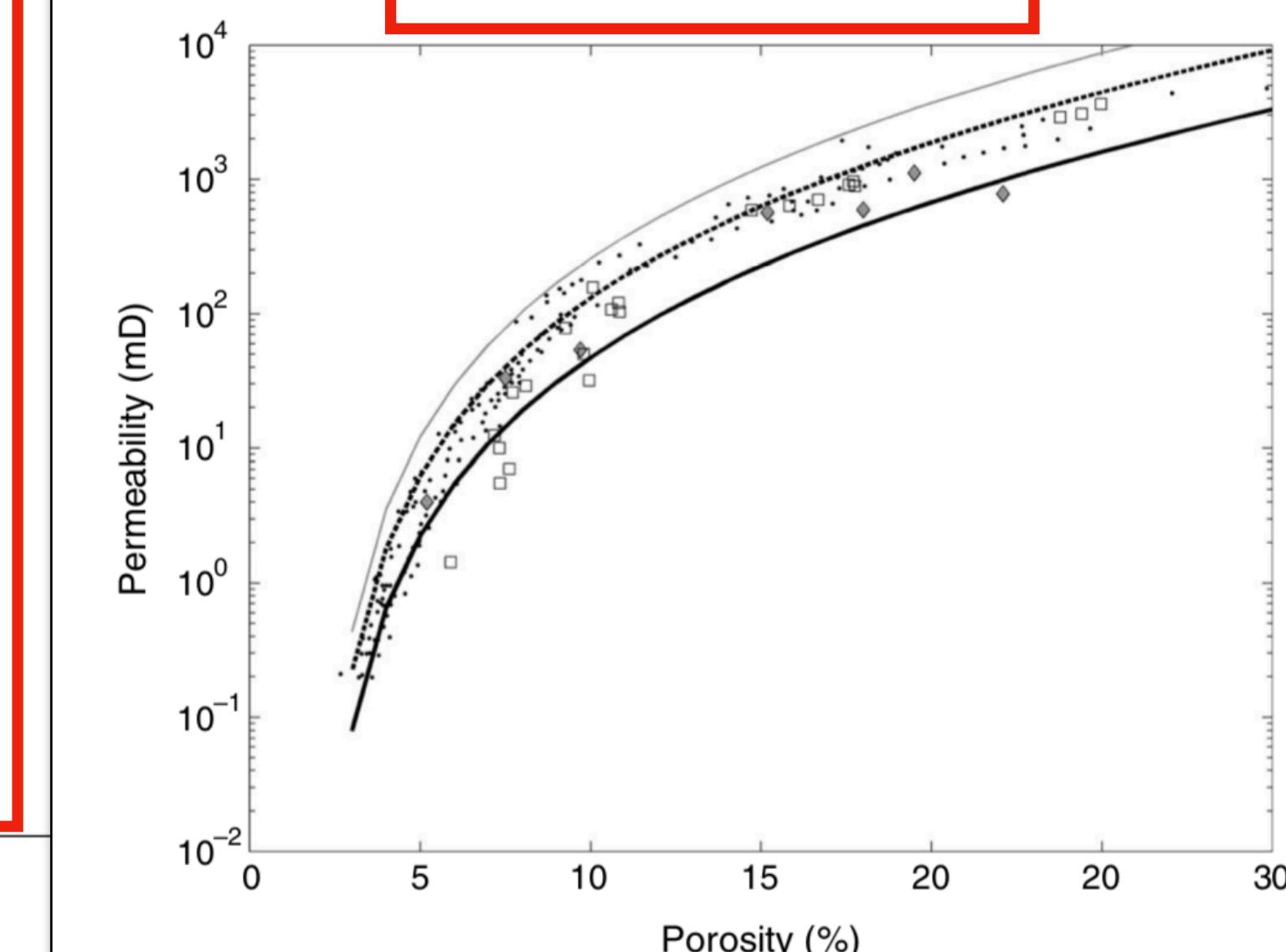
Table 1. Porosity is the helium porosity, permeability R is measured resistivity in ohm·m, F is the normalized (n) of 2, and $R_w = 0.17$ ohm·m, Error (%) is the er

Sample	Porosity	Perm (mD)
A11	0.07	10
A16	0.07	6
A33	0.07	12
A82	0.08	7
A87	0.10	50
A89	0.08	26
A117	0.11	103
B31	0.11	107
B86	0.09	78
B101	0.11	121
B102	0.10	157
B108	0.08	29
F510	0.15	592
GT3	0.17	704
GW18	0.16	637
GW19	0.18	912
GW23	0.18	965
GW28	0.18	896
H27	0.25	3630
H42	0.24	2894
H74	0.24	3079
F410	0.06	1
F570	0.10	32

ABSTRACT

The relations among the resistivity, elastic-wave velocity, porosity, and permeability in Fontainebleau sandstone samples from the Ile de France region, around Paris, France were experimentally revisited. These samples followed a permeability-porosity relation given by Kozeny-Carman's equation. For the resistivity measurements, the samples were partially saturated with brine. Archie's equation was used to estimate resistivity at 100% water saturation, assuming a saturation exponent, $n = 2$. Using

$$m = 0.09 \ln\left(\sqrt{\frac{c\phi^3}{k}}\right) + 1.98. \quad (7)$$



Use Case 2: organic carbon content of Earth's crust



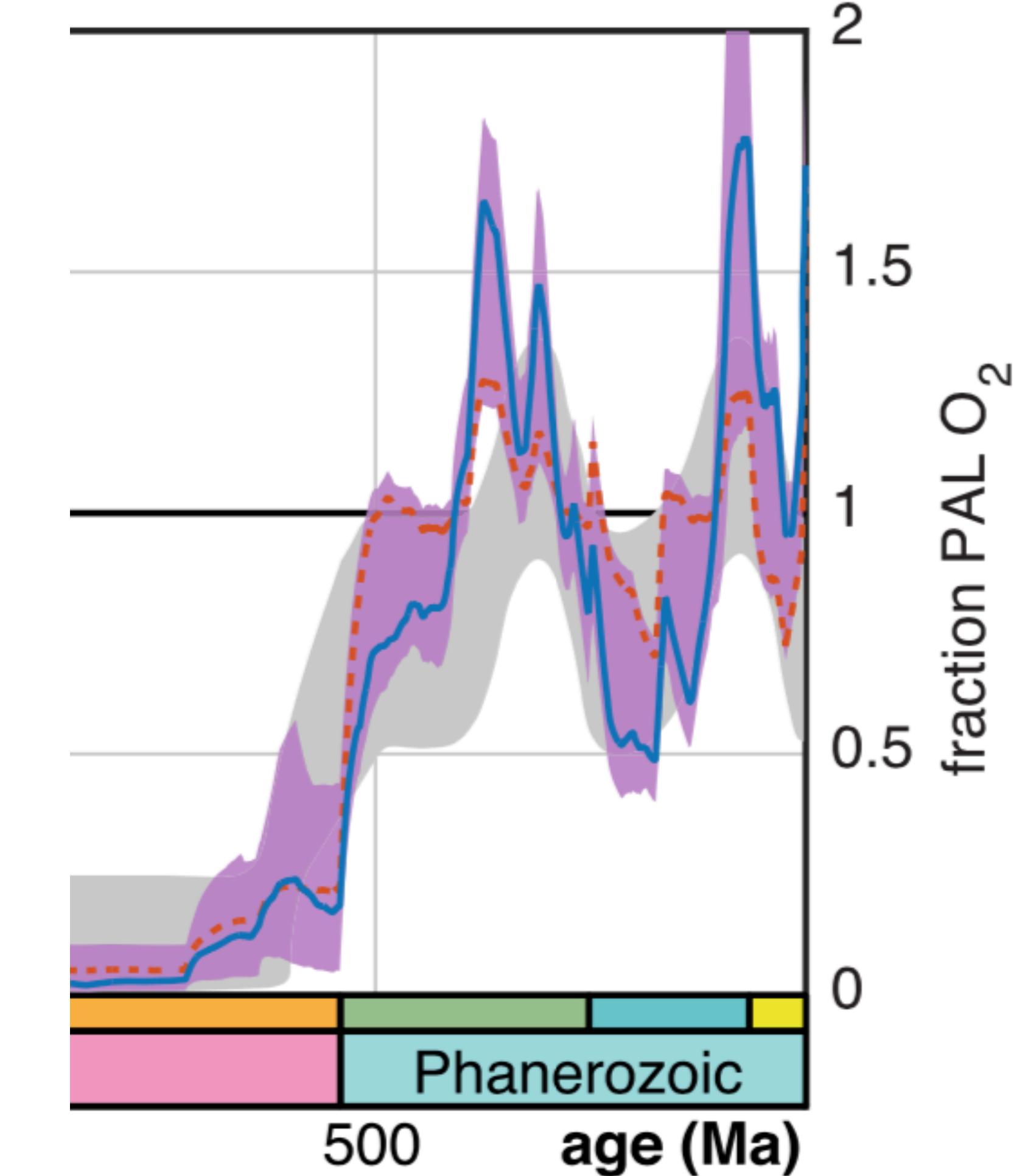
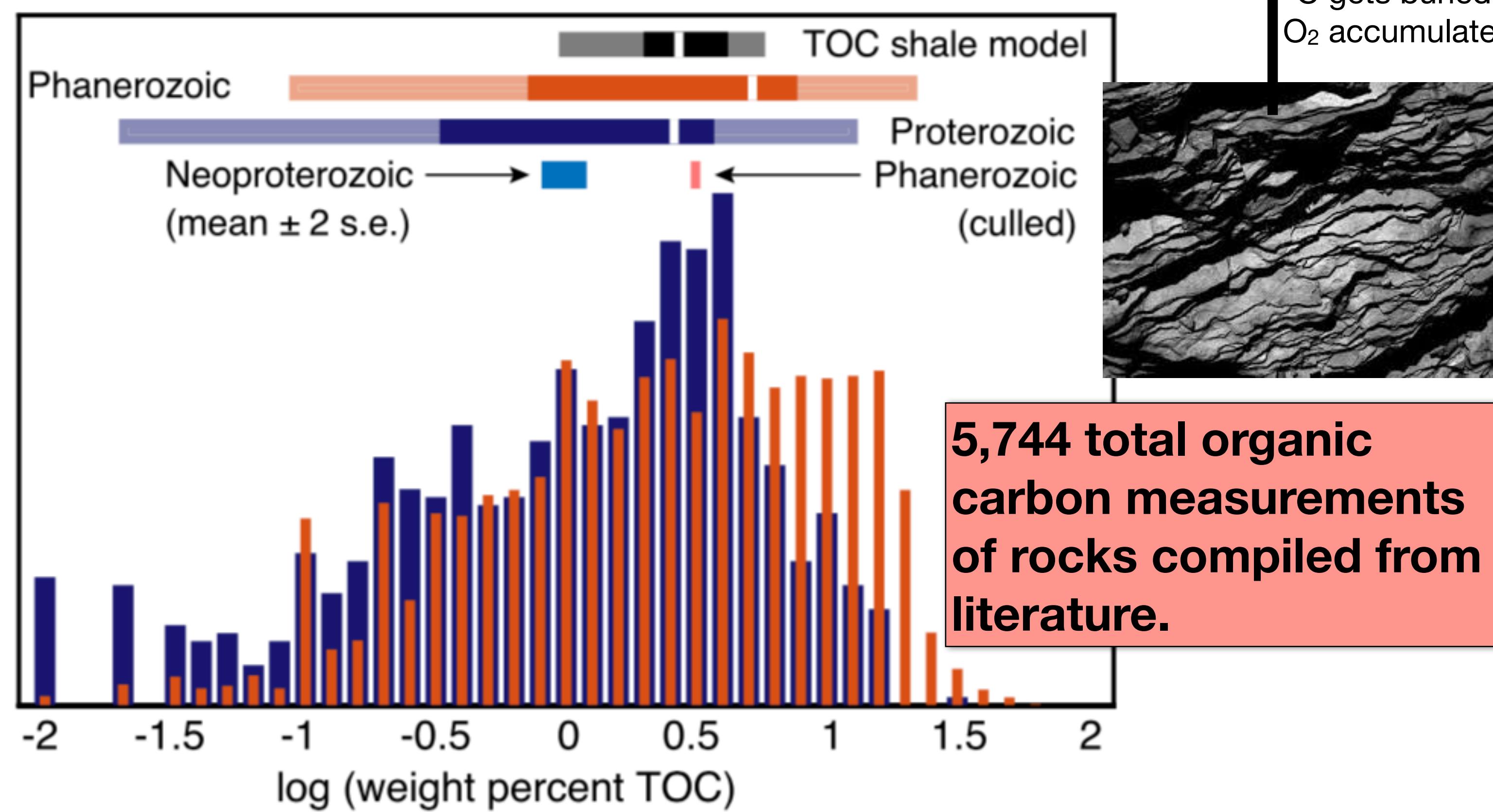
Contents lists available at [ScienceDirect](#)

Earth and Planetary Science Letters

www.elsevier.com/locate/epsl

Atmospheric oxygenation driven by unsteady growth of the continental sedimentary reservoir

Jon M. Husson*, Shanan E. Peters



$$\frac{dM}{dt} = F_{org} - k_1 M - k_2 B_{org} \sqrt{M},$$

Use Case 2: organic carbon content of Earth's crust



Macrostrat



Improve inventories of unconventional fossil fuels, improve models of atmospheric oxygen

Automate location and aggregation of contextualized sample-based measurements of organic carbon in rocks

COSMOS INPUT: vocabulary for specific measurement types and example context, xDD document acquisition/processing pipeline

Total Organic Carbon (TOC), Pyrolysis Data, and Organic Matter Types, Sites 535 (3450 m water depth) and 540 (2926 m water depth)

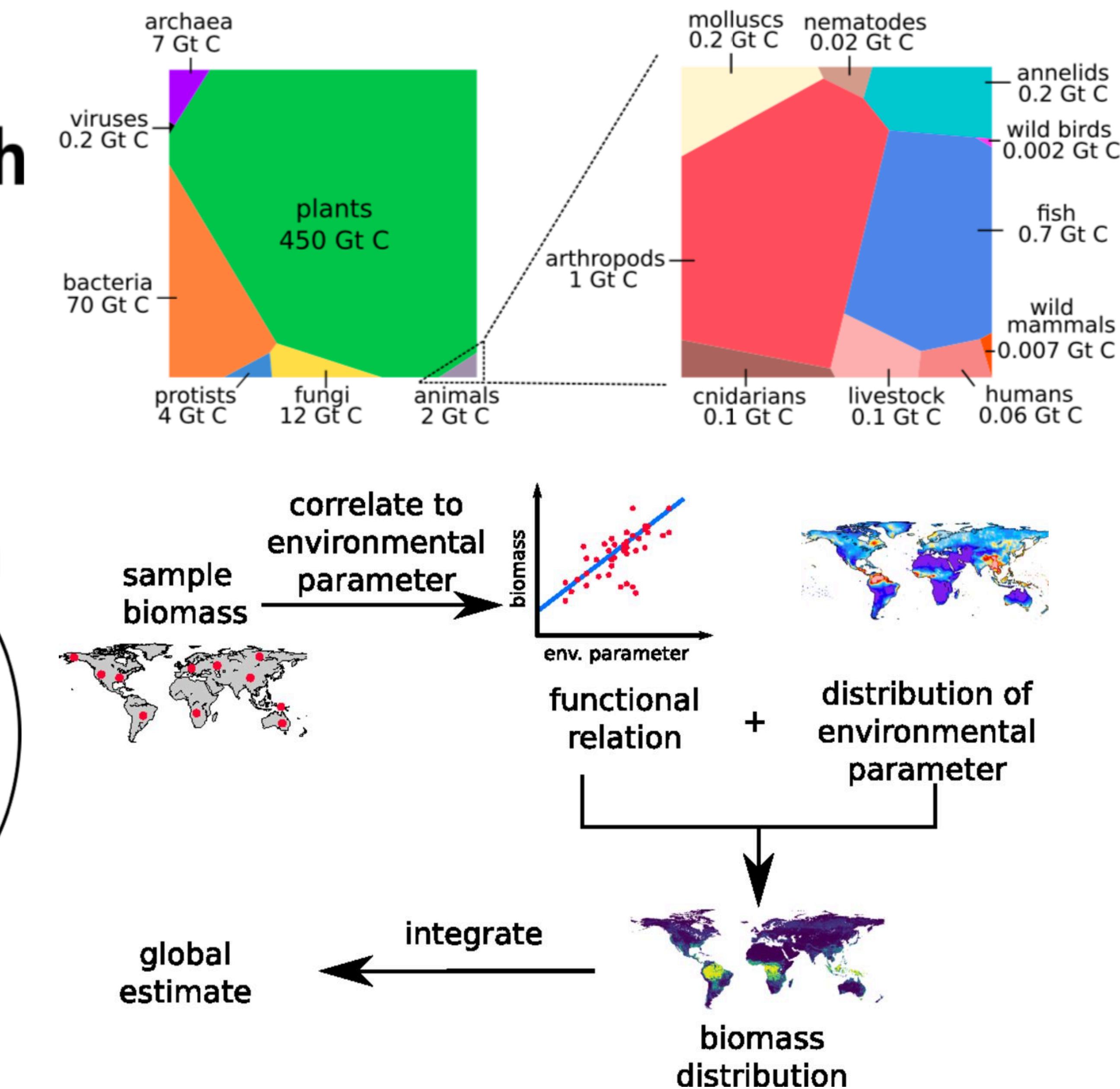
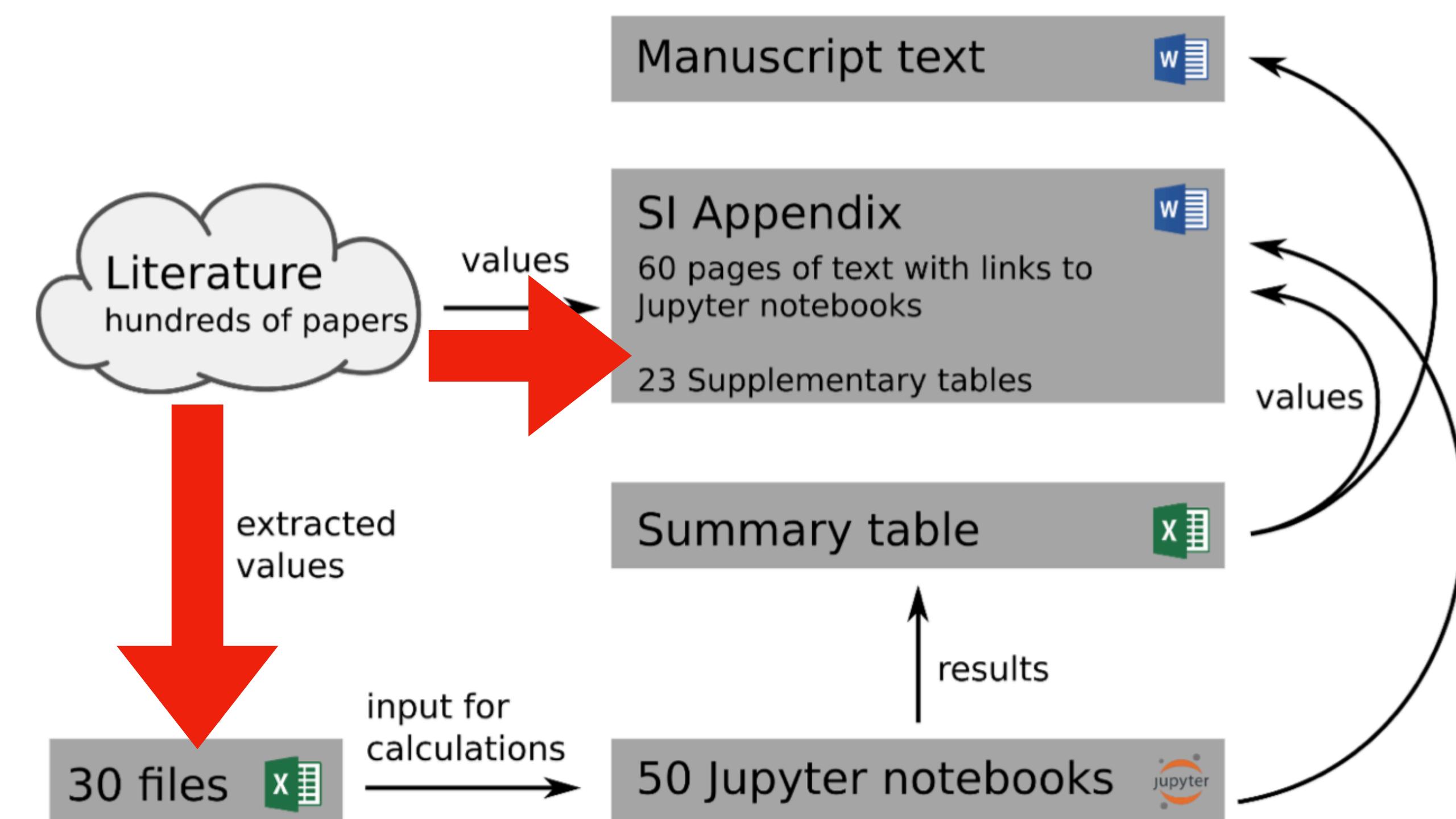
Use Case 3: biomass and its distribution on Earth

PNAS

Proceedings of the
National Academy of Sciences
of the United States of America

The biomass distribution on Earth

Yinon M. Bar-On^a, Rob Phillips^{b,c}, and Ron Milo^{a,1}



Use Case 3: biomass and its distribution on Earth



SCIENCE FOR THE BENEFIT OF HUMANITY



**Automate location
and aggregation of
contextualized
sample-based
measurements of
biomass of specific
organisms**

**Improve biomass
inventories and build
better models**

COSMOS INPUT: vocabulary for specific measurement types, list of target taxonomic names, example context, xDD document acquisition/processing pipeline

Table 1

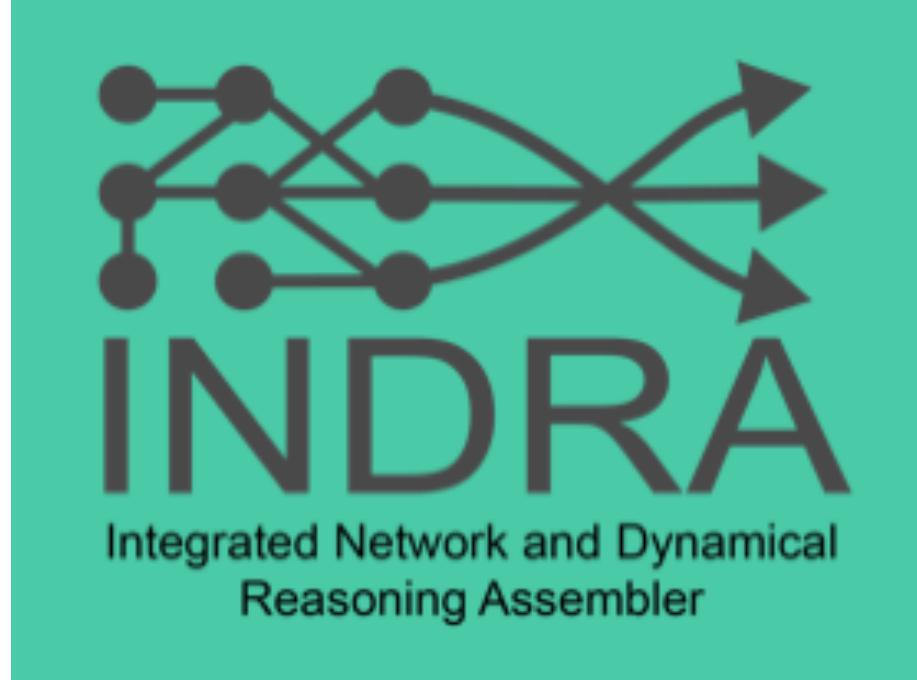
Relative distribution (%) of soil macrofauna biomass in litter and soil layers, total density (ind m^{-2}) and total biomass (g m^{-2}) in a forest (F), 4 years pasture (P4), and abandoned pasture (Pa) in Central Amazon

	F	P4	Pa 95	Pa 97
Relative distribution of macrofaunal biomass in litter (%)	3 (1)	0 (0)	0.5 (0)	0 (0)
Relative distribution of soil macrofaunal biomass in 0–5 cm (%)	13 (7)	6 (4)	27 (16)	4 (3)
Relative distribution of soil macrofaunal biomass in 5–10 cm (%)	3 (2)	19 (11)	44.5 (36.2)	12 (8)
Relative distribution of soil macrofaunal biomass in 10–25 cm (%)	81 (58)	75 (61)	28 (22)	84 (65)
Density of termites (ind m^{-2})	3608 (2113)	941 (689)	2117 (1455)	2266 (1652)
Density of ants (ind m^{-2})	2519 (2010)	322 (145)	2174 (1784)	402 (206)
Density of Oligochaeta (ind m^{-2})	136 (70)	284 (114)	0 (0)	0 (0)
Density of <i>P. corethrurus</i> (ind m^{-2})	0 (0)	0 (0)	390 (224)	136 (87)
Density of Isopoda (ind m^{-2})	35 (27)	5 (3)	38 (29)	110 (86)
Density of Diplopoda (ind m^{-2})	49 (28)	22 (17)	26 (19)	93 (78)
Density of Chilopoda (ind m^{-2})	54 (32)	0 (0)	158 (122)	10 (7)
Biomass of termites (g m^{-2})	5.29 (2.87)	0.61 (0.45)	3.43 (2.02)	4.28 (3.12)
Biomass of ants (g m^{-2})	1.27 (0.98)	0.19 (0.14)	1.57 (1.21)	0.51 (0.30)
Biomass of Oligochaeta (g m^{-2})	44.32 (26.24)	48.91 (29.21)	0 (0)	0 (0)
Biomass of <i>P. corethrurus</i> (g m^{-2})	0 (0)	0 (0)	45.12 (36.14)	7.76 (5.23)
Biomass of Isopoda (g m^{-2})	0.02 (0.01)	0.21 (0.19)	0.05 (0.03)	1.28 (1.02)
Biomass of Diplopoda (g m^{-2})	0.13 (0.06)	0.33 (0.25)	0.22 (0.13)	1.35 (1.11)
Biomass of Chilopoda (g m^{-2})	0.96 (0.46)	0 (0)	0.55 (0.32)	0.07 (0.02)

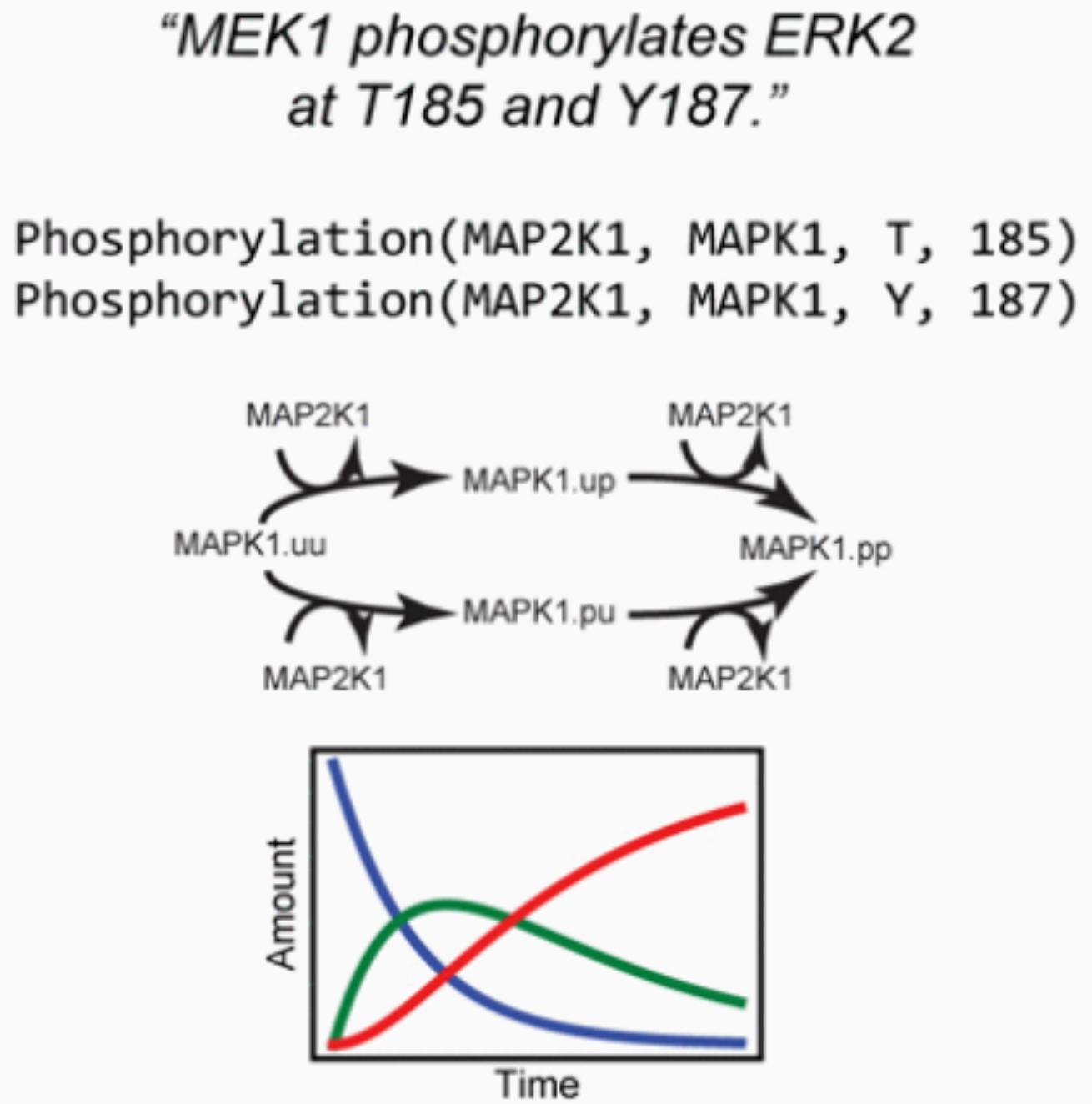
Means of 10 values, standard errors in parentheses.

Use Case 4: INDRA statements

COSMOS INPUT: large, general vocabulary, INDRA code, xDD document acquisition/processing pipeline



```
"title": "Regulated SUMOylation and Ubiquitination of DdMEK1 Is Required for Proper Chemotaxis",
"doi": "10.1016/S1534-5807(02)00186-7",
"coverDate": "June 2002",
"URL": "http://www.sciencedirect.com/science/article/pii/S1534580702001867",
"authors": "Sobko, Alex; Ma, Hui; Firtel, Richard A.",
"hits": 5,
"highlight": [
    ", San Diego 9500 Gilman Drive La Jolla, California 92093 Summary MEK1, which is required for aggregation",
    " to chemoattractant stimulation. SUMOylation is required for MEK1's function and its translocation from the",
    " to the cytosol and cortex, including the leading edge of chemotaxing cells. MEK1 in which the site",
    " of SUMOylation is mutated is retained in the nucleus and does not complement the mek1 null phenotype",
    ". Constitutively active MEK1 is cytosolic and is constitutively SUMOylated, whereas the corresponding"
```



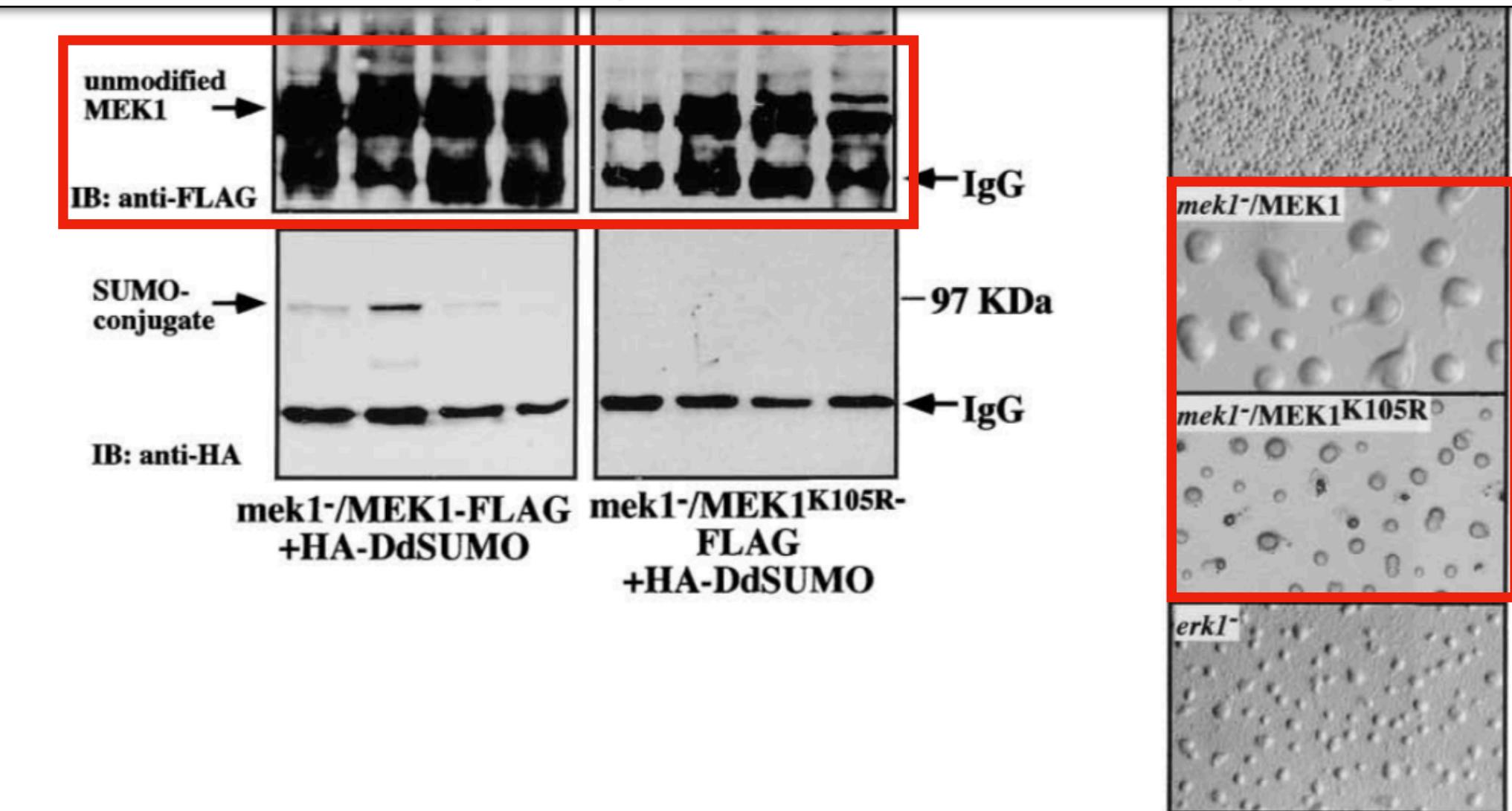
Knowledge

INDRA

Model

API-based access to high-quality full-text mentions

Pipeline to generate INDRA statements



C.

Genotype (Strain)	Speed (μm/min)	Directional change (deg)	Roundness (%)	Directionality
Wild-type (KAx-3)	8.5±0.5	12.5±4.4	43.8±4.2	0.90±0.6
<i>mek1</i> ^{-/-} /MEK1 (ASF1)	8.1±0.4	12.4±0.8	41.3±4.8	0.89±0.05
<i>mip1</i> null (HMF3)	5.7±0.3	41.1±4.5	51.9±0.7	0.59±0.06
<i>mek1</i> ^{-/-} /MEK1K105R	5.0±0.4	42.6±4.4	45.9±6.1	0.66±0.04