



Hypothesis management in support of the e-scientific method

Bernardo Gonçalves
bng@br.ibm.com

Workshop eScience @UFABC – June 22, 2017



IBM: São Paulo's Building in Vila Mariana



The e-Scientific Method

*“Originally, there was just experimental science, and then there was theoretical science, with Kepler’s Laws, Newton’s Laws of Motion, Maxwell’s equations, and so on. Then, for many problems, the **theoretical models** grew too complicated to solve analytically, and **people had to start simulating**. These simulations have carried us through much of the last half of the last century. At this point, these **simulations are generating a whole lot of data**, along with a huge increase in data from the experimental sciences.” — Jim Gray, 2007*

Hypothesis Management



- Scientific research is based on the central idea of a hypothesis, meant to be established or refuted.

Hypothesis Management



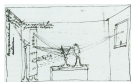
- Scientific research is based on the central idea of a hypothesis, meant to be established or refuted.
- Over time and from multiple sources, we may collect evidence that support their (gray-shaded) truth or falsehood.

Hypothesis Management



- Scientific research is based on the central idea of a **hypothesis**, meant to be established or refuted.
- Over time and from multiple sources, we may collect **evidence** that support their (gray-shaded) truth or falsehood.
- *Hypothesis management* can be, therefore, closely related to the management of **probabilistic** data.

Hypothesis Management



- Scientific research is based on the central idea of a **hypothesis**, meant to be established or refuted.
- Over time and from multiple sources, we may collect **evidence** that support their (gray-shaded) truth or falsehood.
- *Hypothesis management* can be, therefore, closely related to the management of **probabilistic** data.
- ...A '**crucial experiment**' allegedly establishes the truth of one of a set of competing theories (Routledge Encyc. of Philosophy).

1 Research Vision

- Hypothesis Management
- Use Case

2 Probabilistic DB Construction Pipeline

- Hypothesis Encoding
- Probabilistic DB Synthesis
- Prototype System

3 Conclusions

- Takeaways
- Future Work
- Appendix

From Hypotheses to Data

Law of free fall

"If a body falls from rest, then its velocity at any point is proportional to the time it has been falling."

(i)

```
for k = 0:n;
    t = k * dt;
    v = -g*t + v_0;
    s = -(g/2)*t^2 + v_0*t + s_0;
    t_plot(k) = t;
    v_plot(k) = v;
    s_plot(k) = s;
end
```

(iii)

$$a = -g$$

$$v = -g t + v_0$$

$$s = -(g/2) t^2 + v_0 t + s_0$$

(ii)

FALL	t	v	s
	0	0	5000
	1	-32	4984
	2	-64	4936
	3	-96	4856
	4	-128	4744

(iv)

Rival Hypotheses: a Probability Distribution

- Rival hypotheses supposed to explain the same phenomenon.

H₁. Free fall law

$$a = -g$$

$$v = -g t + v_0$$

$$s = -(g/2) t^2 + v_0 t + s_0$$

FALL	ϕ	v	t	v	s
	1	1	0	0	5000.00
	1	1	1	-32	4984.00
	1	1	2	-64	4936.00

$Pr := 0.33$

H₂. Stokes' law

$$a = 0$$

$$v = -\sqrt{gD / 4.6 \times 10^{-4}}$$

$$s = -t \sqrt{gD / 4.6 \times 10^{-4}} + s_0$$

FALL	ϕ	v	t	v	s
	1	2	0	-607.9E-3	5000.00
	1	2	1	-607.9E-3	4997.32
	1	2	2	-607.9E-3	4994.64

$Pr := 0.33$

H₃. Velocity-squared law

$$a = 0$$

$$v = -gD^2 / 3.29 \times 10^{-6}$$

$$s = -(gD^2 / 3.29 \times 10^{-6}) t + s_0$$

FALL	ϕ	v	t	v	s
	1	3	0	-4.17	5000.00
	1	3	1	-4.17	4981.63
	1	3	2	-4.17	4963.26

$Pr := 0.33$

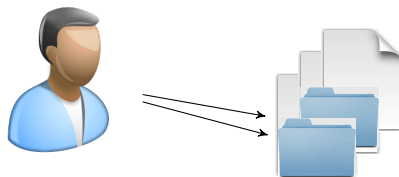
FALL	ϕ	t	s
	1	0	5000
	1	1	4979.3
	1	2	4932.6

Concrete Use Scenario in Computational Science

Example 1.

Bob is a computational scientist who is playing with a number of models and different parameter settings to see which one gives a best fit to his observation samples.

Each run constitutes a specific model instantiation that is associated with a unique file ('big table'). In the end of the day his resulting datasets are spread over many files and folders. □



Beyond Files: a Big Table Database

- User's **default** choice: struggle with the files to find relevant data.

“There is life beyond files.” (Jim Gray)

FALL	tid	t	g	v_0	s_0	a	v	s
	1	0	32	0	5000	-32	0	5000
	1	1	32	0	5000	-32	-32	4984
	1	2	32	0	5000	-32	-64	4936

	2	0	32.2	0	5000	-32.2	0	5000
	2	1	32.2	0	5000	-32.2	-32.2	4983.9
	2	2	32.2	0	5000	-32.2	-64.4	4935.6

Do Better: a Hypothesis Database

$$a = -g$$

$$v = -g t + v_0$$

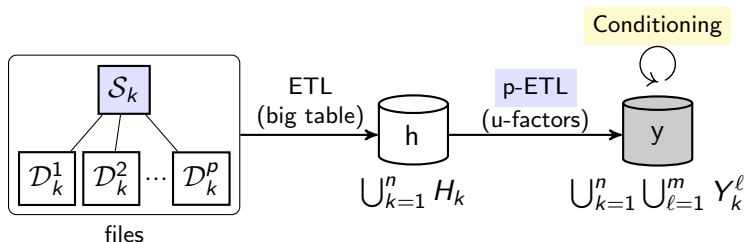
$$s = -(g/2) t^2 + v_0 t + s_0$$

FALL	tid	t	g	v ₀	s ₀	a	v	s
	1	0	32	0	5000	-32	0	5000
	1	1	32	0	5000	-32	-32	4984
	1	2	32	0	5000	-32	-64	4936

	2	0	32.2	0	5000	-32.2	0	5000
	2	1	32.2	0	5000	-32.2	-32.2	4983.9
	2	2	32.2	0	5000	-32.2	-64.4	4935.6

- Predictive structure: strong **correlations** from the math models (!).
- Project-level **standardization**: pick a favorite MathML editor to report and manage model equations declaratively.

Design-by-Synthesis Pipeline



Technical challenges:

- ① **Encoding**: math equations \rightarrow structural eqs. \rightarrow functional deps.;
- ② **Causal reasoning**: inferring the causal ordering and u-factors;
- ③ **Probabilistic DB synthesis**: normalization based on the u-factors;
- ④ **Conditioning**: probability distribution update in face of evidence.

- 1 Research Vision
 - Hypothesis Management
 - Use Case
- 2 Probabilistic DB Construction Pipeline
 - Hypothesis Encoding
 - Probabilistic DB Synthesis
 - Prototype System
- 3 Conclusions
 - Takeaways
 - Future Work
 - Appendix

Given set \mathcal{E} of equations over set \mathcal{V} of variables... (1)

Law of free fall:

$\mathcal{H} = \{$

$$g = 32,$$

$$v_0 = 0,$$

$$s_0 = 5000,$$

$$a = -g,$$

$$v = -g t + v_0,$$

$$s = -(g/2)t^2 + v_0 t + s_0 \}.$$

$\xrightarrow{h-encode}$

$\Sigma = \{$

$$\phi \rightarrow g,$$

$$\phi \rightarrow v_0,$$

$$\phi \rightarrow s_0,$$

$$g \ v \rightarrow a,$$

$$g \ v_0 \ t \ v \rightarrow v,$$

$$g \ v_0 \ s_0 \ t \ v \rightarrow s \}.$$

Given set \mathcal{E} of equations over set \mathcal{V} of variables... (2)

Law of free fall:

$\mathcal{S} = \{$

$$f_1(g) = 0,$$

$$f_2(v_0) = 0,$$

$$f_3(s_0) = 0,$$

$$f_4(a, g) = 0,$$

$$f_5(v, g, t, v_0) = 0,$$

$$f_6(s, g, t, v_0, s_0) = 0 \}.$$

$\xrightarrow{h\text{-encode}}$

$\Sigma = \{$

$$\phi \rightarrow g,$$

$$\phi \rightarrow v_0,$$

$$\phi \rightarrow s_0,$$

$$g \ v \rightarrow a,$$

$$g \ v_0 \ t \ v \rightarrow v,$$

$$g \ v_0 \ s_0 \ t \ v \rightarrow s \}.$$

Causal Ordering Algorithm (AI Literature)

- ① Identify '**minimal substructures**' at step k ;
- ② Reduce the matrix by eliminating them;
- ③ Call step $k+1$ **recursively**.

	x_1	x_2	x_3	x_4	x_5	x_6	x_7
f_1	1	0	0	0	0	0	0
f_2	0	1	0	0	0	0	0
f_3	0	0	1	0	0	0	0
f_4	1	1	1	1	1	0	0
f_5	1	0	1	1	1	0	0
f_6	0	0	0	1	0	1	0
f_7	0	0	0	0	1	0	1



	x_1	x_2	x_3	x_4	x_5	x_6	x_7
f_1	1	0	0	0	0	0	0
f_2	0	1	0	0	0	0	0
f_3	0	0	1	0	0	0	0
f_4	1	1	1	1	1	0	0
f_5	1	0	1	1	1	0	0
f_6	0	0	0	1	0	1	0
f_7	0	0	0	0	1	0	1

Equivalent to Finding a Biclique $K_{m,n}$ in a Bipartite Graph

$$S = \{$$

$$f_1(x_1) = 0,$$

$$f_2(x_2) = 0,$$

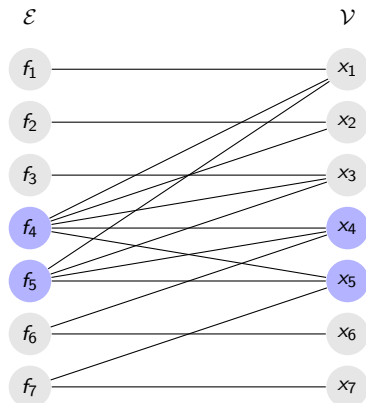
$$f_3(x_3) = 0,$$

$$f_4(x_1, x_2, x_3, x_4, x_5) = 0,$$

$$f_5(x_1, x_3, x_4, x_5) = 0,$$

$$f_6(x_4, x_6) = 0,$$

$$f_7(x_5, x_7) = 0 \}.$$

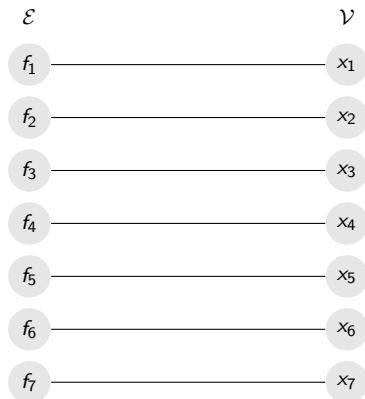


Theorem 1.

Let $S(\mathcal{E}, \mathcal{V})$ be a complete structure. The extraction of its causal ordering by $\text{COA}(S)$ tries to solve an NP-Hard problem.

Easier: Complete Matching in a Bipartite Graph

- Hopcroft-Karp algorithm to solve it in $O(\sqrt{|\mathcal{E}|} |\mathcal{S}|)$;



Proposition 2.

Let $\mathcal{S}(\mathcal{E}, \mathcal{V})$ be a structure. Then a total causal mapping $\varphi: \mathcal{E} \rightarrow \mathcal{V}$ over \mathcal{S} exists iff \mathcal{S} is complete.

Provably Correct Approach to Hypothesis Encoding

$$C_\varphi = \{ (x_a, x_b) \mid \text{there exists } f \in \mathcal{E} \text{ such that } \varphi(f) = x_b \\ \text{and } x_a \in \text{Vars}(f) \}$$

Proposition 1.

Let $\mathcal{S}(\mathcal{E}, \mathcal{V})$ be a structure, and $\varphi_1: \mathcal{E} \rightarrow \mathcal{V}$ and $\varphi_2: \mathcal{E} \rightarrow \mathcal{V}$ be any two total causal mappings over \mathcal{S} . Then $C_{\varphi_1}^+ = C_{\varphi_2}^+$.

Causal Ordering: Sub-Quadratic Complexity on $|\mathcal{S}|$

Corollary 1.

Let $\mathcal{S}(\mathcal{E}, \mathcal{V})$ be a complete structure. Then a **total causal mapping** $\varphi: \mathcal{E} \rightarrow \mathcal{V}$ over \mathcal{S} can be found by $\text{TCM}(\mathcal{S})$ in time that is **bounded by $O(\sqrt{|\mathcal{E}|} |\mathcal{S}|)$** .



Artificial Intelligence

Volume 238, September 2016, Pages 154-165



A note on the complexity of the causal ordering problem

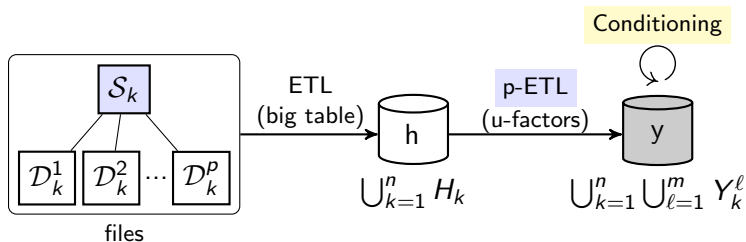
Bernardo Gonçalves^a, Fabio Porto^b

Show more

<https://doi.org/10.1016/j.artint.2016.06.004>

[Get rights and content](#)

Coming Back from Detail



Technical challenges:

- ① **Encoding**: math equations \rightarrow structural eqs. \rightarrow functional deps.;
- ② **Causal reasoning**: inferring the causal ordering and u-factors;
- ③ **Probabilistic DB synthesis**: normalization based on the u-factors;
- ④ **Conditioning**: probability distribution update in face of evidence.

Bob's Big Table (before our help)

H ₃	tid	ϕ	v	t	x_0	b	p	y_0	d	r	x	y
	1	1	3	0	30	.5	.02	4	.75	.02	30	4
	1	1	3	...	30	.5	.02	4	.75	.02
	2	1	3	0	30	.5	.018	4	.75	.023	30	4
	2	1	3	...	30	.5	.018	4	.75	.023
	3	1	3	0	30	.4	.02	4	.8	.02	30	4
	3	1	3	...	30	.4	.02	4	.8	.02
	4	1	3	0	30	.4	.018	4	.8	.023	30	4
	4	1	3	...	30	.4	.018	4	.8	.023
	5	1	3	0	30	.397	.02	4	.786	.02	30	4
	5	1	3	...	30	.397	.02	4	.786	.02
	6	1	3	0	30	.397	.018	4	.786	.023	30	4
	6	1	3	5	30	.397	.018	4	.786	.023	50.1	62.9
	6	1	3	10	30	.397	.018	4	.786	.023	13.8	8.65
	6	1	3	15	30	.397	.018	4	.786	.023	79.3	8.23
	6	1	3	20	30	.397	.018	4	.786	.023	12.6	30.7
	6	1	3	...	30	.397	.018	4	.786	.023

Synthesized Tables: Ready for Predictive Analysis

Y_0	$V \mapsto D$	ϕ	v
	$x_1 \mapsto 1$	1	1
	$x_1 \mapsto 2$	1	2
	$x_1 \mapsto 3$	1	3

Y_3^1	$V \mapsto D$	ϕ	x_0
	$x_2 \mapsto 1$	1	30

Y_3^2	$V \mapsto D$	ϕ	b
	$x_3 \mapsto 1$	1	.5
	$x_3 \mapsto 2$	1	.4
	$x_3 \mapsto 3$	1	.397

Y_3^3	$V \mapsto D$	ϕ	p
	$x_4 \mapsto 1$	1	.020
	$x_4 \mapsto 2$	1	.018

Y_3^4	$V_1 \mapsto D_1$	$V_2 \mapsto D_2$	$V_3 \mapsto D_3$	$V_4 \mapsto D_4$	ϕ	v	t	y	x
	$x_1 \mapsto 3$	$x_2 \mapsto 1$	$x_3 \mapsto 1$	$x_4 \mapsto 1$	1	3	1900	4	30
	$x_1 \mapsto 3$	$x_2 \mapsto 1$	$x_3 \mapsto 1$	$x_4 \mapsto 1$	1	3
	1	3
	$x_1 \mapsto 3$	$x_2 \mapsto 1$	$x_3 \mapsto 3$	$x_4 \mapsto 2$	1	3	1900	4	30
	$x_1 \mapsto 3$	$x_2 \mapsto 1$	$x_3 \mapsto 3$	$x_4 \mapsto 2$	1	3	1901	4.12	41.5
	$x_1 \mapsto 3$	$x_2 \mapsto 1$	$x_3 \mapsto 3$	$x_4 \mapsto 2$	1	3	1902	5.78	56.7
	$x_1 \mapsto 3$	$x_2 \mapsto 1$	$x_3 \mapsto 3$	$x_4 \mapsto 2$	1	3	1903	11.7	72.8
	$x_1 \mapsto 3$	$x_2 \mapsto 1$	$x_3 \mapsto 3$	$x_4 \mapsto 2$	1	3	1904	31.1	75.9
	$x_1 \mapsto 3$	$x_2 \mapsto 1$	$x_3 \mapsto 3$	$x_4 \mapsto 2$	1	3

Querying Rival Predictions with Probabilities

Y_3^4	$V_1 \mapsto D_1$	$V_2 \mapsto D_2$	$V_3 \mapsto D_3$	$V_4 \mapsto D_4$	ϕ	v	t	y	x
	$x_1 \mapsto 3$	$x_2 \mapsto 1$	$x_3 \mapsto 1$	$x_4 \mapsto 1$	2	3	1900	4	30

	$x_1 \mapsto 3$	$x_2 \mapsto 1$	$x_3 \mapsto 3$	$x_4 \mapsto 2$	2	3	1904	..	75.92

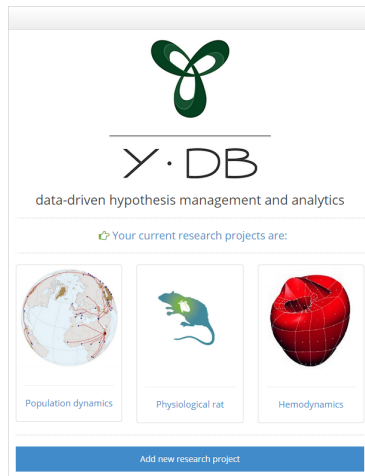
W	$V \mapsto D$	Pr
	$x_1 \mapsto 2$.33
	$x_1 \mapsto 3$.33
	$x_1 \mapsto 3$.33
	$x_2 \mapsto 1$	1
	$x_3 \mapsto 1$.33
	$x_3 \mapsto 2$.33
	$x_3 \mapsto 3$.33
	$x_4 \mapsto 1$.5
	$x_4 \mapsto 2$.5

Operation **conf()**;

$\theta = \{x_1 \mapsto 3, x_2 \mapsto 1, x_3 \mapsto 3, x_4 \mapsto 2\},$

$Pr = .33 * 1 * .33 * .5 \approx .055$

Prototype System (1)



Prototype System (2)

Phenomenon data definition

Phenomenon id

Research

Description

Upload dataset (CSV format)
 Lynx_Hare.csv

Loading observations

70%

Observable
Year
Lynx

Prototype System (3)

Hyphotesis data definition

Hypothesis id

Name

Upload structure (XML format)

Lotka_Volterra.xml

Hypothesis structure: 100%

Phenomenon

Map symbols

Variable	Observable
<input type="text" value="t"/>	<input type="text" value="Year"/>
<input type="text" value="x"/>	<input type="text" value="Lynx"/>

Hypothesis trial datasets (MAT format)

10 files

Prototype System (4)

Hypothesis management

Hypothesis

Lotka-Volterra model

Phenomenon

Lynx population in Hudson's Bay, Canada, from 1900 to 1920.

Simulation trial

6

Key 1 Key 2

« 1 2 3 4 5 6 7 8 9 10 ... 20 »

t	y	x
1904.0	31.1083920070696	75.9196961932191
1904.1	34.1895828779035	74.4878105315043
1904.2	37.4356490187431	72.6675156977604
1904.3	40.8008312965646	70.4705874385644
1904.4	44.2262696820135	67.9230251780413
1904.5	47.6417755935588	65.0649920502872
1904.6	50.9691309633260	61.9493743700121
1904.7	54.1267031673551	58.6389813583626

Prototype System (5)

Hypothesis Analytics			
Phenomenon			
Lynx population in Hudson's Bay, Canada, from 1900 to 1920.			
Observations Predictions			
« 1 2 3 »			
	Year	Lynx	Hare
<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/>	1900	30	4
<input checked="" type="checkbox"/>	1901	47.2	6.1
<input checked="" type="checkbox"/>	1902	70.2	9.8
<input checked="" type="checkbox"/>	1903	77.4	35.2
<input checked="" type="checkbox"/>	1904	36.3	59.4
<input checked="" type="checkbox"/>	1905	20.6	41.7
<input checked="" type="checkbox"/>	1906	18.1	19
<input checked="" type="checkbox"/>	1907	21.4	13
<input checked="" type="checkbox"/>	1908	22	8.3

Prototype System (6)

Hypothesis Analytics

Phenomenon
Lynx population in Hudson's Bay, Canada, from 1900 to 1920. ▼

Observations Predictions

« 1 2 3 4 5 6 7 8 9 10 ... 20 »

upsilon	tid	Year	Lynx	conf
3	2	1904	65.060410460081	0.183505
3	6	1904	75.919696193219	0.179993
3	4	1904	77.459735769215	0.175992
3	1	1904	89.592307430943	0.131452
3	5	1904	88.321831841064	0.127000
3	3	1904	90.083803232660	0.124023
1	1	1904	16.487212706992	0.047211
2	2	1904	77.822475573932	0.017372
2	1	1904	79.812581025093	0.013234
1	2	1904	18.221188003898	0.000220

- 1 Research Vision
 - Hypothesis Management
 - Use Case
- 2 Probabilistic DB Construction Pipeline
 - Hypothesis Encoding
 - Probabilistic DB Synthesis
 - Prototype System
- 3 **Conclusions**
 - **Takeaways**
 - Future Work
 - Appendix

Takeaways

We have seen:

- The automatic construction of a probabilistic DB (out of math equations and datasets) to support the analysis of crucial exps .
- This is hypothesis management (as data management and analytics) in support of the e-scientific method .

(Papers: *PVLDB*'14, *IEEE Computing in Science & Eng.*'15, *Artif. Intell.*'16)

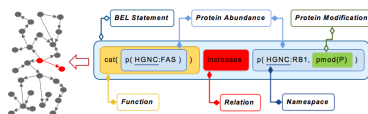
- 1 Research Vision
 - Hypothesis Management
 - Use Case
- 2 Probabilistic DB Construction Pipeline
 - Hypothesis Encoding
 - Probabilistic DB Synthesis
 - Prototype System
- 3 Conclusions
 - Takeaways
 - **Future Work**
 - Appendix

Future Work: in the field of Bioinformatics

- Recommendation of crucial experiments;
 - Find 'rival' (structurally similar) math models from a repository;
 - Example: BioModels (EMBL-EBI), with 1.6K+ models stored.

Future Work: in the field of Bioinformatics

- **Recommendation** of crucial experiments;
 - Find 'rival' (structurally similar) math models from a repository;
 - Example: BioModels (EMBL-EBI), with 1.6K+ models stored.
- **Refutation** attempts (sense of Karl Popper);
 - Look for negative claims in the literature.
 - Example: Causal Biological Networks, with 120+ models stored (each has hundreds of hypothetical claims).



SET Citation = "PubMed", "Regulation of Rb and cdk2 by signal transduction cascades: divergent effects of JNK1 and p38 kinases." , "EMBO J, 1999 Mar 15;18(6):1559-70." , "1007592"

SET Evidence = "The stimulation of Jurkat cells is known to induce p38 kinase and we find a pronounced increase in Rb phosphorylation within 30 min of Fas stimulation"

SET Tissue = "Jurkat cells"

Thank you.

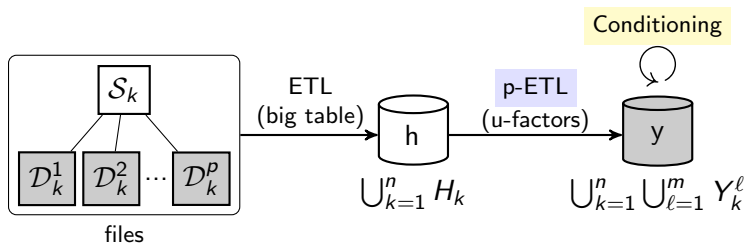
Acknowledgements



Questions?

Bernardo Gonçalves
bng@br.ibm.com

Design-by-Synthesis Pipeline



- Technical challenges:

- ① **Encoding:** math equations \rightarrow structural eqs. \rightarrow functional deps.;
- ② **Causal reasoning:** inferring the causal ordering and u-factors;
- ③ **Probabilistic DB synthesis:** normalization based on the u-factors;
- ④ **Conditioning:** probability distribution update in face of evidence.

The Folding Σ^{\rightarrow} of Σ

Acyclic pseudo-transitive reasoning

Algorithm 1 Folding of an fd set.

```
1: procedure FOLDING( $\Sigma$ : fd set)
Require:  $\Sigma$  given encodes complete structure  $\mathcal{S}$ 
Ensure: Returns fd set  $\Sigma^{\rightarrow}$ , the folding of  $\Sigma$ 
2:    $\Sigma^{\rightarrow} \leftarrow \emptyset$ 
3:   for all  $\langle X, A \rangle \in \Sigma$  do
4:      $Z \leftarrow \text{AFolding}(\Sigma, A)$ 
5:      $\Sigma^{\rightarrow} \leftarrow \Sigma^{\rightarrow} \cup \langle Z, A \rangle$ 
6:   return  $\Sigma^{\rightarrow}$ 
```

The Folding $\Sigma^{\text{q}\rightarrow}$ of Σ

Acyclic pseudo-transitive reasoning

$$\Sigma = \{ \phi \rightarrow x_0,$$

$$\phi \rightarrow b,$$

$$\phi \rightarrow p,$$

$$\phi \rightarrow y_0,$$

$$\phi \rightarrow d,$$

$$\phi \rightarrow r,$$

$$x_0 b p \text{ } t v y \rightarrow x,$$

$$y_0 d r \text{ } t v x \rightarrow y \}.$$

folding
→

$$\Upsilon(\Sigma)^{\text{q}\rightarrow} = \{$$

$$x_0 b p \text{ } t v y \text{ } y_0 d r \rightarrow x,$$

$$y_0 d r \text{ } t v x \text{ } x_0 b p \rightarrow y \}.$$

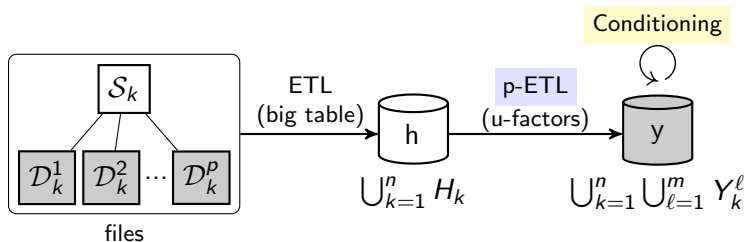
FD set Σ_{89} (a real Physiology Model that “fits the screen”)

$$\Sigma_{89} = \{ \begin{array}{l} \phi \rightarrow C1a \ C1p \ C2a \ C2p \ C3a \ Cglobal \ Cmyo \\ \quad Dp100 \ Pc \ t_delta \ t_max \ t_min \ taua \ taud, \\ \quad \phi \ t \rightarrow DelP, \\ C1a \ C1p \ C2a \ C2p \ C3a \ Cglobal \ Cmyo \ Dp100 \ Pc \ v \rightarrow Dc, \\ \quad Dc \ Pc \ v \rightarrow Tc, \\ Cglobal \ Cmyo \ Dc \ Pc \ v \rightarrow Ac, \\ Dc \ v \rightarrow D_t_min, \\ Ac \ v \rightarrow A_t_min, \\ DelP \ Pc \ v \rightarrow P, \\ D \ P \ v \rightarrow T, \\ A \ C1a \ C1p \ C2a \ C2p \ C3a \ Dp100 \ P \ T \ v \rightarrow Ttarget, \\ \quad Cglobal \ Cmyo \ D \ P \ v \rightarrow Atarget, \\ D_t_min \ Dc \ T \ Tc \ Ttarget \ t \ taud \ v \rightarrow D, \\ A_t_min \ Atarget \ t \ taua \ v \rightarrow A \}. \end{array}$$

Its Folding $\Sigma_{89}^{\rightarrow}$

$$\Upsilon(\Sigma_{89})^{\rightarrow} = \{ \begin{array}{l} C1a \ C2a \ \phi \ v \ \rightarrow \ A_t_min \ Ac \ D_t_min \ Dc \ Tc, \\ C1a \ DelP \ \phi \ v \ \rightarrow \ P, \\ C1a \ C2a \ DelP \ \phi \ t \ v \ T \rightarrow \ A \ Atarget \ D \ Ttarget \end{array} \}.$$

Design-by-Synthesis Pipeline



- Technical challenges:

- ① **Encoding:** math equations \rightarrow structural eqs. \rightarrow functional deps.;
- ② **Causal reasoning:** inferring the causal ordering and u-factors;
- ③ **Probabilistic DB synthesis:** normalization based on the u-factors;
- ④ **Conditioning:** probability distribution update in face of evidence.

Defining the Theoretical U-factor

$$Y_0 := \pi_{\phi, v}(\text{repair-key}_{\phi}(H_0)).$$

H_0	ϕ	v
	1	1
	1	2
	1	3

Y_0	$V \mapsto D$	ϕ	v
	$x_1 \mapsto 1$	1	1
	$x_1 \mapsto 2$	1	2
	$x_1 \mapsto 3$	1	3

W	$V \mapsto D$	Pr
	$x_1 \mapsto 1$.33
	$x_1 \mapsto 2$.33
	$x_1 \mapsto 3$.33

U-factor Learning

Discovery of contingent functional dependencies

- **'Input'** relation of the Lotka-Volterra hypothesis. Observe:
 - **Multiplicity** of parameter values;
 - **Correlations** between parameter values.

H_3^1	tid	ϕ	x_0	b	p	y_0	d	r
	1	1	30	.5	.020	4	.75	.020
	2	1	30	.5	.018	4	.75	.023
	3	1	30	.4	.020	4	.8	.020
	4	1	30	.4	.018	4	.8	.023
	5	1	30	.397	.020	4	.786	.020
	6	1	30	.397	.018	4	.786	.023

$$\Omega = \{ \begin{array}{l} \phi x_0 \rightarrow y_0, \\ \phi b \rightarrow d, \\ \phi p \rightarrow r \end{array} \}.$$

U-factorization

Defining the Empirical U-factors

$$Y_k^i := \pi_{\phi A G} (\text{repair-key}_{\phi @ \text{count}} (\gamma_{\phi, A, G, \text{count}(\ast)}(H_k))).$$

Y_3^1	$V \mapsto D$	ϕ	x_0	y_0
	$x_2 \mapsto 1$	1	30	4

Y_3^2	$V \mapsto D$	ϕ	b	p
	$x_3 \mapsto 1$	1	.5	.5
	$x_3 \mapsto 2$	1	.4	.8
	$x_3 \mapsto 3$	1	.397	.786

Y_3^3	$V \mapsto D$	ϕ	p	r
	$x_4 \mapsto 1$	1	.020	.020
	$x_4 \mapsto 2$	1	.018	.023

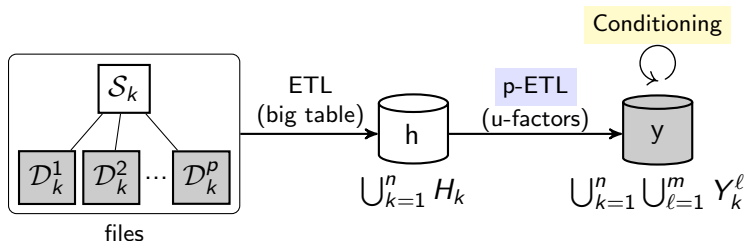
W	$V \mapsto D$	Pr

	$x_2 \mapsto 1$	1
	$x_3 \mapsto 1$.33
	$x_3 \mapsto 2$.33
	$x_3 \mapsto 3$.33
	$x_4 \mapsto 1$.5
	$x_4 \mapsto 2$.5

Design-Theoretic Properties

- Desirable properties for probability update are ensured;
 - **Claim-centered** decomposition.
 - Theorem 6: **BCNF** w.r.t. causal dependencies.
 - **Correctness** of uncertainty decomposition.
 - Theorem 7: **Lossless join** w.r.t. causal dependencies.

Design-by-Synthesis Pipeline



• Technical challenges:

- ① **Encoding:** math equations \rightarrow structural eqs. \rightarrow functional deps.;
- ② **Causal reasoning:** inferring the causal ordering and u-factors;
- ③ **Probabilistic DB synthesis:** normalization based on the u-factors;
- ④ **Conditioning:** probability distribution update in face of evidence.

Systematic Application of Bayesian Inference

- ① User selection of the **observation** sample;
- ② System selection of the competing **prediction** samples;
- ③ Bayesian **inference** ;
- ④ Probability distribution **update** ;

Bayes' Rule

- Normal density **likelihood** function:

$$f(y | \mu_k) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y-\mu_k)^2} \quad (1)$$

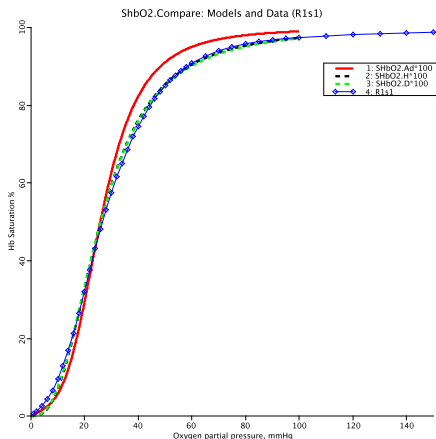
- Bayes'** rule:

$$p(\mu_k | y_1, \dots, y_n) = \frac{\prod_{j=1}^n f(y_j | \mu_{kj}) p(\mu_k)}{\sum_{i=1}^m \prod_{j=1}^n f(y_j | \mu_{ij}) p(\mu_i)} \quad (2)$$

where y_1, \dots, y_n is the observation sample, μ_k is prediction k and σ is the standard deviation parameter.

Probability Update in Face of Evidence

STUDY	ϕ	v	pO2	SHbO2	Prior	Posterior
	1	32	100	9.72764121981342E-1	.333	.335441
	1	28	100	9.74346796798538E-1	.333	.335398
	1	31	100	9.90781330988763E-1	.333	.329161



Viewpoint: Why Hypothesis Management?



"Numerical simulations and 'big data' are essential in modern science, but they do not alone yield understanding. Building a massive database to feed simulations without corrective loops between hypotheses and experimental tests seems, at best, a waste of time and money." **Nature** 513, Sept 2014.