



HUMAN COMPUTATION FOR MATERIALS DISCOVERY



Ronan Le Bras Richard Bernstein Carla P. Gomes Bart Selman R. Bruce van Dover Computer Science Computer Science Computer Science Computer Science Materials Science/Physics

December 7, 2012

NIPS Workshop - HCSCS



cfci

Cornell Fuel Cell Institute

Mission: develop **new materials** for **fuel cells**.



Figure 1. Fuel cell schematic. Source: Annual Reveiws of Energy and the Environment. http://energy.annualreviews.org/ cgi/content/full/24/1/281

- An **Electrocatalyst** must:
- 1) Be electronically conducting
- 2) Facilitate both reactions
- **Platinum** is the best known metal to fulfill that role, but:
- 1) The reaction rate is still considered slow (causing **energy loss**)
- 2) Platinum is fairly **costly**, **intolerant** to fuel **contaminants**, and has a **short lifetime**.

Goal: Find an intermetallic compound that is a better catalyst than Pt.





Recipe for finding alternatives to Platinum

- 1) In a vacuum chamber, place a silicon wafer.
- 2) Add three metals.
- 3) Mix until smooth, using three sputter guns.
- *4)* Bake for 2 hours at 650°C



- *Deliberately* inhomogeneous composition on Si wafer
- Atoms are intimately mixed







Figure 1: Phase regions of Ta-Rh-Pd



Figure 2: Fluorescence activity of Ta-Rh-Pd





Identifying crystal structure using **X-Ray Diffraction** at CHESS

• The XRD pattern **characterizes** the underlying **lattice** at a given point on the silicon wafer.









Scattered

details

beam

Bragg's law: $\begin{array}{c} \underset{heam}{\overset{\text{lncident}}{\longrightarrow}} \\ \lambda = 3.0 \\ & \end{array} \\ & d = 3.0 \\ & A \end{array}$

Atomic planes







Atomic planes





Outline



✓ Motivation

- Physical Characteristics
- Problem Definition
- Satisfiability Modulo Theory Approach
- UDiscoverIt
- Empirical Validation
- Conclusions and Future Work







































Additional Physical characteristics:

- Phase Connectivity
- Mixtures of < 3 pure phases
- Peaks shift by \leq 15% within a region
 - Continuous and Monotonic
- Noisy detection of peaks





- ✓ Motivation
- ✓ Physical Characteristics
- Problem Definition
- Satisfiability Modulo Theory Approach
- UDiscoverIt
- Empirical Validation
- Conclusions and Future Work





• Input:

• A list of points on the silicon wafer



- A real vector D_i per vertex v_i (diffraction patterns)
- *K* = user specified number of pure phases
- Goal: a basis of K vectors for

X vectors for

$$D_1 \quad D_2 \quad \dots \quad D_N$$

 $D_i = a_{il}B_1 + \dots + a_{ik}B_K$







• Non-negative basis vectors and coefficients

 $\boldsymbol{B_i \geq 0}$, $a_{ij} \geq 0$

• At most M (=3) non-zero coefficients per point

 $|\{j \mid a_{ij} > 0\}| \le M$

• Basis patterns appear in **contiguous** locations on silicon wafer Build a graph *G* of the points on the silicon wafer $V_1 = V_2 = V_3$ The subgraph induced by $|\{i \mid a_{ij} > 0\}|$ is connected *G*



 V_{\varDelta}





• Shifts coefficients are **bounded**, **continuous** and **monotonic**







- ✓ Motivation
- ✓ Physical Characteristics
- ✓ Problem Definition
- Satisfiability Modulo Theory Approach
- UDiscoverIt
- Empirical Validation
- Conclusions and Future Work





• Initial graph G and number K of basis patterns







- Initial graph *G* and number *K* of basis patterns
- *Peak detection* to extract a set of peaks P_i for each diffraction pattern D_i





ICS is

- Initial graph G and number K of basis patterns
- *Peak detection* to extract a set of peaks P_i for each diffraction pattern D_i
- Real variable e_{jk} for the location of peak k in basis B_j







- Initial graph *G* and number *K* of basis patterns
- *Peak detection* to extract a set of peaks P_i for each diffraction pattern D_i
- Real variable e_{ik} for the location of peak k in basis B_i
- Real variable s_{ij} for the **shift** coefficient of basis B_j in point P_i





- An observed peak p is "*explained*" if there exists s_{ij} , e_{jk} s.t. $|p (s_{ij} + e_{jk})| \le \varepsilon$
- Every observed peak must be *explained*







- An observed peak p is "*explained*" if there exists s_{ij} , e_{jk} s.t. $|p (s_{ij} + e_{jk})| \le \varepsilon$
- Every observed peak must be *explained*
- Some peaks might be missing (unobserved)
- Bound the number of missing peaks $\leq T$
- Minimization by (binary) search on T







- Linear phase usage constraints (up to *M* basis patterns per point)
- Linear constraints for shift monotonicity and continuity ($s_{ij} \le s_{lm}$)
- Lazy connectivity: add a cut if current solution is not connected

If disconnected regions explained with phase 1



Then Phase 1 must appear in at least one of these points

- Symmetry breaking:
 - Renaming of pure phases
 - Ordering of the peak locations e_{ik} (per basis pattern)

> Quantifier-free linear arithmetic





	Time (s)					
System	P	L^*	K	#var	#cst	
A/B/C	36	8	4	408	2095	3502
A/B/C	60	8	4	624	3369	17345
Al/Li/Fe	15	6	6	267	1009	79
Al/Li/Fe	28	6	6	436	1864	346
Al/Li/Fe	28	8	6	490	2131	10076
Al/Li/Fe	28	10	6	526	2309	28170
Al/Li/Fe	45	7	6	693	3281	18882
Al/Li/Fe	45	8	6	711	3410	46816

Table: Runtime (seconds) of the SMT solver. P is the number of sample points, L^* is the average number of peaks per phase, K is the number of basis patterns, #var is the number of variables and #cst is the number of constraints.





- Can human input be used to significantly boost the performance of combinatorial reasoning and optimization methods?.
- Can human input provide useful global guidance to the solver, by identifying the setting of so-called **backdoor variables** in the SMT model?





- ✓ Motivation
- ✓ Physical Characteristics
- ✓ Problem Definition
- ✓ Satisfiability Modulo Theory Approach
- UDiscoverIt
- Empirical Validation
- Conclusions and Future Work





- Graphical User Interface for providing user input to an SMT solver
 - Visualization of the X-ray diffraction data
 - User provides **partial assignments** of the variables of the SMT formulation $\begin{bmatrix} r & 3 & 1 = 1 \end{bmatrix}$

```
 \begin{array}{c} & & \\ (\text{let } ((?x3935 \ (+ \ (+ \ ?x3913 \ (\text{ite } (\text{and } r_3_1 \ (\text{not } \$x3920)) \ 1 \ 0)) \ (\text{ite } (\text{and } r_3_1 \ (\text{not } \$x3931)) \ 1 \ 0)))) \\ & & \\ (\text{let } ((?x3957 \ (+ \ (+ \ ?x3935 \ (\text{ite } (\text{and } r_3_1 \ (\text{not } \$x3942)) \ 1 \ 0)) \ (\text{ite } (\text{and } r_3_1 \ (\text{not } \$x3953)) \ 1 \ 0)))) \\ & & \\ (\text{let } ((\$x994 \ (\text{and } (<= \ (+ \ e_{-}1_{-}9 \ S_{-}2_{-}1) \ 6442) \ (<= \ 6432 \ (+ \ e_{-}1_{-}9 \ S_{-}2_{-}1))))) \\ & & \\ (\text{let } ((\$x991 \ (+ \ e_{-}1_{-}9 \ S_{-}2_{-}1))) \\ & & \\ (\text{let } ((\$x1080 \ (<= \ 7372 \ ?x991)))) \\ & & \\ (\text{let } ((\$x1079 \ (<= \ 7391 \ 7382)))) \\ & & \\ (\text{let } ((\$x1081 \ (\text{and } \$x1079 \ \$x1080)))) \\ & & \\ (\text{let } ((\$x1168 \ (\text{and } (<= \ ?x991 \ 7747) \ (<= \ 7737 \ ?x991))))) \end{array}
```





- Graphical User Interface for providing user input to an SMT solver
 - Visualization of the X-ray diffraction data
 - User provides **partial assignments** of the variables of the SMT formulation
 - **Representation** inspired from how **materials scientists analyze** the data and address this problem









- Graphical User Interface for providing user input to an SMT solver
 - Visualization of the X-ray diffraction data
 - User provides **partial assignments** of the variables of the SMT formulation
 - **Representation** inspired from how **materials scientists analyze** the data and address this problem
 - No knowledge requirement about the underlying phase structure, and only very limited knowledge of diffraction methods.



UDiscoverIt

















- ✓ Motivation
- ✓ Physical Characteristics
- ✓ Problem Definition
- \checkmark Satisfiability Modulo Theory Approach
- ✓ UDiscoverIt
- Empirical Validation
- Conclusions and Future Work





Dataset						Time w/o	Time w/	# assigned
System	P	L^*	K	#var	#cst	user input (s)	user input (s)	var. by user
A/B/C	36	8	4	408	2095	3502	150	19 (4.6%)
A/B/C	60	8	4	624	3369	17345	261	18 (2.9%)
Al/Li/Fe	15	6	6	267	1009	79	27	6 (2.2%)
Al/Li/Fe	28	6	6	436	1864	346	83	12 (2.7%)
Al/Li/Fe	28	8	6	490	2131	10076	435	26 (5.3%)
Al/Li/Fe	28	10	6	526	2309	28170	188	23 (4.3%)
Al/Li/Fe	45	7	6	693	3281	18882	105	28 (4.0%)
Al/Li/Fe	45	8	6	711	3410	46816	74	30 (4.2%)

Table: Runtime (seconds) of the SMT solver with and without user input. P is the number of sample points, L^* is the average number of peaks per phase, K is the number of basis patterns, #var is the number of variables and #cst is the number of constraints.





- With **limited effort** and **input**, a user can provide insightful information about the structure of the problem, and dramatically **speed up** the performance of the SMT solver.
- The approach leverages the **complementary strength** of **human input**, providing **global insights** about problem structure, and the power of **combinatorial solvers** to exploit **complex local constraints**.





- Aggregating input from multiple users
- Guiding the search as a variable/value ordering heuristics, as opposed to pre-assignments of variables
- Providing the user with explanation and feedback about inconsistencies
- Adapting the GUI to a more 'FoldIt' spirit
- Adapting this method to Mechanical Turk



Future Work



Artificial Ar	anical turk		Your Account	HITS	Qualifications	398,185 HITs	5
, internet of a	and an		All HITS HITS Ava	ailable To You	HITs Assigned 1	o You	
	Find HITE	Containing			that nav at h	for which	you are qualified
	THE PARTY	contoining			that pay at a		
er: 00:00:00 of	60 minutes		Want to work on Accept HI	this HIT? Want to	o see other HITs? Skip HIT		Total Earned: Unavail Total HITs Submitted: 0
dentify patterns of I	lines in an image		<u></u>				
Requester: UDisco	overIt					Reward: \$0.20 per H	T HITs Available: 50 Duration: 60 minute
Quantications Rec	Juirea: None						
dontific a l	Dattann of Va	utical Lines					
dentify a l	rattern of ver	rtical Lines					
nd the most dis	tinguishable pattern	of vertical lines/blobs	intersecting the targe	t row. Decide	which of the v	ertical lines/blobs are def	initely part of this pattern, possibly part
this pattern, an ategory selected	la separate from this	s pattern, or that you a	are uncertain about (o	overlapping, fai	nt, obscured, e	etc) and mark them by ci	icking on them with the appropriate
this is your fire	t took in this series a	lease read the detailed	d instructions and look	at the example	ec hefore starti	ng the task using the fol	owing buttons. Defer back to the
arifications if u	nclear situations arise		a menucuone and loop	c at the example	es belore starti	ing the task, using the for	lowing buttons. Refer back to the
Instructions	Easier Example	Harder Example	Current HIT				
uniont LUT							
			-				
Definitely Primary Pat	in Pro tern Prima	obably in ary Pattern	Don't Know	Secondary	Pattern	Definitely in Secondary Pattern	
			9	6			
		2 3 4 5	67891011	12 13 14 15	16 17	Hinh	
			67891011	12 13 14 15		High	
				12 13 14 15	16 17	High	
				12 13 14 15		High	
				12 13 14 15		High	
				00		High	
				•••		High	
				0.0		High	
				1 12 13 14 15		High	
				12 13 14 15		High	
				12 13 14 15		High	
				1 12 13 14 15		High	
				1/2 13 14 15		High	
				1 12 13 14 15		High	
				12 13 14 15			
				12 13 14 15			
				12 13 14 15			

Want to work on this HIT? Want to see other HITs?







THANK YOU!

Extra slides









- Motivation
- Problem Definition (Part I)
- Prior Work: Non-negative Matrix Factorization
- Problem Definition (Part II)
- Our Work: Satisfiability Modulo Theories Approach
- Conclusion and Future work



Guidelines





ICS

Non-negative Matrix Factorization [Long et al., 2009]

Advantages: scales up very well, accurately solves simple systems *Drawbacks*: overlooks critical physical behavior, making the results physically meaningless for more complex systems.







- Motivation
- Problem Definition (Part I)
- Prior Work: Non-negative Matrix Factorization
- Problem Definition (Part II)
- Our Work: Satisfiability Modulo Theories Approach
- Conclusion and Future work



Experimental Results



•Illustration on Al-Li-Fe system



Runtime



# Points	Unknown Phases	Arithmetic + Z3 (s)	Set-based + CPLEX (s)
10	3	8	0.5
	6	12	Timeout
15	3	13	0.5
	6	20	Timeout
18	3	29	384.8
	6	125	Timeout
29	3	78	276
	6	186	Timeout
45	6	518	Timeout

Z3 scales to realistic sized problems!

Arithmetic encoding translated to CP and MIP:

- MIP is appealing because it can optimize the objective
- They don't scale \rightarrow SMT solving strategy



Precision/Recall



Li Li Li Li Li Li Li Li	AI	Li Basis Patterns α β Fe γ Li δ ξ	Fround SMT Truth Results
	Size	Precision	Recall
Al Fe Al Fe	10	95.8	100
	15	96.6	100
Recovers around truth	18	97.2	96.6
.	29	96.1	92.8
	45	95.8	91.6





- Remove some peaks to simulate experimental noise
- Size = 15 points



Solutions are still accurate. Runtime increases approx linearly.



Previous Work 1: Cluster Analysis [Long et al., 2007]





PCA – 3 dimensional approx

Hierarchical Agglomerative Clustering

Drawback: Requires sampling of pure phases, detects phase regions (not phases), overlooks peak shifts, may violate physical constraints (phase continuity, etc.).







Drawback: Overlooks peak shifts (linear combination only), may violate physical constraints (phase continuity, etc.).





Parameters

- Number of pure phases K, tolerance ε
- Key components
 - Variables peak positions per base
 - Shifts per point
 - Point *p* is explained by base *k*





- New arithmetic-based encoding:
 - Real variables e_{ii} for the peak locations in each B_i
 - Real variables for the shift coefficients s_{ij} (per base, per point)
 - An observed peak *p* is explained if $|p-s_{ij} e_{ij}| \le \varepsilon$ (Match the height of the peaks)
 - Bound the number of missing peaks $\leq T$

