

# Constraint Reasoning and Kernel Clustering for Pattern Decomposition With Scaling



#### **Ronan LeBras**

Theodoros Damoulas Ashish Sabharwal Carla P. Gomes John M. Gregoire Bruce van Dover Computer Science Computer Science Computer Science Computer Science Materials Science / Physics Materials Science / Physics

Sept 15, 2011

CP'11



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



ICS B

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

- XRD pattern characterizes the underlying crystal fairly well
- **Expensive** experimentations: Bruce van Dover's research team has access to the facility **one week every year**.









































#### **Additional Physical characteristics:**

- Peaks shift by  $\leq 15\%$  within a region
- Phase Connectivity
- Mixtures of ≤ 3 phases
- Small peaks might be discriminative
- Peak locations matter, more than peak intensities









Figure 1: Phase regions of Ta-Rh-Pd

Figure 2: Fluorescence activity of Ta-Rh-Pd



### Outline



- Motivation
- Problem Definition
  - $\checkmark$  Abstraction
  - ✓ Hardness
- CP Model
- Kernel-based Clustering
- Bridging CP and Machine learning
- Conclusion and Future work





Input

- $-G = (V, E) \text{ be an undirected graph with } V = \{v_1, \dots, v_N\},$ 
  - $-\mathcal{P} = \{P_1, \ldots, P_N\}$  be a collection of N patterns over a finite set  $S \subseteq \mathbb{Q}^+$ ,
  - $-M \leq K \leq N$  be positive integers, and  $\delta \geq 1$  be a rational.





Input

- $-G = (V, E) be an undirected graph with V = \{v_1, \dots, v_N\},$ 
  - $-\mathcal{P} = \{P_1, \ldots, P_N\}$  be a collection of N patterns over a finite set  $S \subseteq \mathbb{Q}^+$ ,
  - $-M \leq K \leq N$  be positive integers, and  $\delta \geq 1$  be a rational.







Input -G = (V, E) be an undirected graph with  $V = \{v_1, \ldots, v_N\}$ ,  $-\mathcal{P} = \{P_1, \ldots, P_N\}$  be a collection of N patterns over a finite set  $S \subseteq \mathbb{Q}^+$ ,  $-M \leq K \leq N$  be positive integers, and  $\delta \geq 1$  be a rational.







Input -G = (V, E) be an undirected graph with  $V = \{v_1, \ldots, v_N\}$ ,  $-\mathcal{P} = \{P_1, \ldots, P_N\}$  be a collection of N patterns over a finite set  $S \subseteq \mathbb{Q}^+$ ,  $-M \leq K \leq N$  be positive integers, and  $\delta \geq 1$  be a rational.



 $M = K = 2, \delta = 1.5$ 





- Input -G = (V, E) be an undirected graph with  $V = \{v_1, \ldots, v_N\}$ ,  $-\mathcal{P} = \{P_1, \ldots, P_N\}$  be a collection of N patterns over a finite set  $S \subseteq \mathbb{Q}^+$ ,  $-M \leq K \leq N$  be positive integers, and  $\delta \geq 1$  be a rational.
- **Output** Determine whether there exists a collection  $\mathcal{B}$  of K basis patterns over S and scaling factors  $s_{ik} \in \{0\} \cup [1/\delta, \delta]$  such that:

(a) 
$$\forall i: P_i = \bigcup_{k=1}^{K} s_{ik} B_k$$
  
(b)  $\forall i: |\{k \mid s_{ik} > 0\}| \le M$   
(c)  $\forall k:$  the subgraph of G induced by  $V_k = \{v_i \in V \mid s_{ik} > 0\}$  is connected  
 $B_1 \longrightarrow P_1 \longrightarrow P_1$ 







- Input -G = (V, E) be an undirected graph with  $V = \{v_1, \ldots, v_N\}$ ,  $-\mathcal{P} = \{P_1, \ldots, P_N\}$  be a collection of N patterns over a finite set  $S \subseteq \mathbb{Q}^+$ ,  $-M \leq K \leq N$  be positive integers, and  $\delta \geq 1$  be a rational.
- Output Determine whether there exists a collection  $\mathcal{B}$  of K basis patterns over S and scaling factors  $s_{ik} \in \{0\} \cup [1/\delta, \delta]$  such that:



M = K = 2,



### **Problem Hardness**







### Outline



- Motivation
- Problem Definition
- CP Model
  - ✓ Model

#### ✓ Experimental Results

- Kernel-based Clustering
- Bridging CP and Machine learning
- Conclusion and Future work





Va	riables	Descriptio	on				-
$p_{ki}$	Normalizi	Normalizing element for phase $k$ in pattern $P_i$					
$a_{ki}$		Whether	r phase k is present in pattern $P_i$				
$q_k$ Set of no			ormalized peak locations of phase $B_k$				
							-
$(a_{ki} = 0)$	$\Leftrightarrow_{K} (p)$	$k_i = 0$		$\forall \ 1 \leq k$	$\leq K, 1 \leq$	$i \leq n$	(1)
$1 \leq$	$\sum^{n} a_{si}$	$\leq M$		$\forall 1$	$\leq i \leq n$		(2)
$(p_{ki} = j) \Rightarrow ($	$q_k \subseteq r_{ij}$	)	$\forall \ 1 \leq k$	$\leq K, 1 \leq$	$i \leq n, 1$	$\leq j \leq  P_i $	(3)
$(p_{ki} = j \land \sum^{K}$	$a_{si} = 1$	$) \Rightarrow (r_{ij} \subseteq$	$q_k)$				
s=1			$\forall \ 1 \leq$	$k \le K, 1$	$\leq i \leq n,$	$1 \le j \le  $	$P_i $ (4)



# CP Model (continued)



$$\begin{aligned} p_{ki} &= j \wedge p_{k'i} = j' \wedge \sum_{s=1}^{K} a_{si} = 2) \\ &\Rightarrow \left( member(r_{ij}[j''], q_k) \lor member(r_{ij'}[j''], q_{k'}) \right) \\ &\forall 1 \le k, k' \le K, 1 \le i \le n, 1 \le j, j', j'' \le |P_i| \quad (5) \end{aligned}$$

$$(p_{ki} = j) \Rightarrow (p_{ki'} \ne j') \qquad \forall 1 \le k \le K, (i, j, i', j') \in \Phi \qquad (6) \\ \Phi = \{ (i, j, i', j') \mid \frac{P_i[j]}{P_i'[j']} < 1/\delta \lor \frac{P_i[j]}{P_i'[j']} > \delta, i < i' \} \end{aligned}$$

$$basisPatternConnectivity(\{a_{ki} | 1 \le i \le n\}) \qquad \forall 1 \le k \le K \qquad (7) \end{aligned}$$

*Advantage:* Captures physical properties and relies on peak location rather than height. *Drawback:* Does not scale to realistic instances; poor propagation if experimental noise.







For realistic instances, K' = 6 and  $N \approx 218...$ 



### Outline



- Motivation
- Problem Definition
- CP Model
- Kernel-based Clustering
- Bridging CP and Machine learning
- Conclusion and Future work



# Kernel-based Clustering



Set of features:  $X = \alpha$ 



#### Similarity matrices:

 $[X.X^T]$ 





#### Method: K-means on Dynamic Time-Warping kernel

**Goal:** Select groups of samples that belong to the same **phase region** to feed the CP model, in order to extract the underlying phases of these sub-problems.





Goal: a robust, *physically meaningful*, scalable, automated solution method that combines:





### Bridging Constraint Reasoning and Machine Learning: Overview of the Methodology





#### **Experimental Validation**







#### Hybrid approach to clustering under constraints

- More robust than data-driven "global" ML approaches
- More scalable than a pure CP model "locally" enforcing constraints

#### An exciting application in close collaboration with physicists

- Best inference out of expensive experiments
- Towards the design of better fuel cell technology





#### **Ongoing work:**

**Spatial Clustering,** to further enhance cluster quality **Bayesian Approach**, to better exploit prior knowledge about local smoothness and available inorganic libraries

#### **Active learning:**

Where to sample next, assuming we can interfere with the sampling process?

When to stop sampling if sufficient information has been obtained?

#### **Correlating catalytic properties across many thin-films:**

In order to understand the underlying **physical mechanism** of **catalysis** and to find promising intermettalic compounds







# Thank you!



# Extra slides







#### Example on Al-Li-Fe diagram:





# Applications with similar structure





#### Flight Calls / Bird conservation

Identifying bird population from sound recordings at night.

Analogy: basis pattern = species samples = recordings physical constraints = spatial constraints, species and season specificities...

#### Fire Detection

Detecting/Locating fires.

Analogy: basis pattern = warmth sources samples = temperature recordings physical constraints = gradient of temperatures, material properties...









<sup>(</sup>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.).

