

# CROWDSOURCING INSIGHTS INTO PROBLEM STRUCTURE FOR SCIENTIFIC DISCOVERY

---



**Bart Selman**

**Carla P. Gomes**

**Ronan Le Bras (now @ AI2)**

**Stefano Ermon (now @ Stanford)**

# Motivation: Automated Scientific Discovery

---



*Can we use modern automated reasoning tools for scientific discovery?*

Science and scientific discovery are arguably among the most impressive demonstrations of human intelligence.

We will consider two areas:

- 1) Finite Mathematics**
- 2) Materials Science**

*In the early days of computing, there was significant optimism with the discovery of first-order theorem proving procedures.*

**But, the raw size of search spaces for interesting math proved to be a daunting challenge. In a sense, the representation language (first-order logic) is too general and unconstrained to effectively explore.**

**Still, some success stories of automated mathematical discoveries, involving “finite mathematics.” (Objects of interest are finite.)**

**We only consider cases of true discovery of previously unknown results (i.e. no “re-discovery”).**

# Discoveries in Discrete Math



**A) 1976 --- Four Color Thm.** Four colors suffices to color a planar map. (Haken and Appel) Problem was open since **1852**. Essentially solved using an exhaustive “proof by cases.” (two thousand special sub-maps were analyzed. (90% human / 10% machine)

Proof is now accepted and has been replicated. Considered somewhat unsatisfying / “inelegant” from a math perspective, because the proof does not provide further insights into why the thm is true. Also, how do we know the computer proof is correct? (Do we?)

**B) 1996 --- Robbins algebras are Boolean algebras.** Concerns sets of axioms. (McCune) Open since **1933**. Proof strategy: clever syntactic re-write of axioms. Several months of cpu time to find the re-write steps. “Semi-human” readable. (30% human / 70% machine)

**C) 2014 --- Erdos Discrepancy Problem.** Property of sequences of +1/-1.  
E.g. +1 +1 -1 +1 -1 -1 +1 -1 ...

Does the sum (absolute value) over subsequences of this sequence stay within a certain bound, C?

Erdos conjectured (**1932**): If sequence gets long enough “discrepancy” exceeds any value of C. Lisitsa and Konev (2014) showed every sequence of **1161 or longer** exceeds bound of **2** (i.e.  $C = 2$ ). *But, remarkably, they also found a +1/-1 sequence of **1160** elements that stayed within bound of 2.* (10% human / 90% machine)

Obtained with Boolean **Satisfiability** solver (20 hrs). Proof trace of **13 gigabytes**. Longest proof in history. Proof **verified** by proof checker. (Contrast two previous results.)

*Each result represents a clear new mathematical advance. But, unfortunately, approaches do not provide much insight into the underlying problem domain.*

For example, can we understand **how** certain complex objects (e.g. 1160 long  $+1/-1$  sequence or coloring of sub-map) can be constructed directly using a human-understandable strategy?

# Example Domain:

## The *Spatially-balanced Latin square (SBLS)* problem

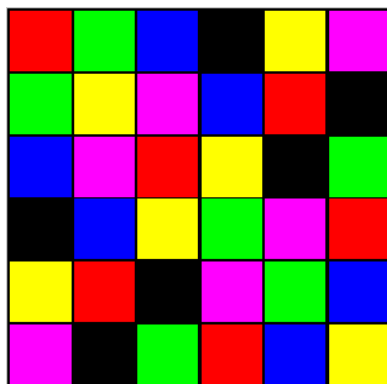


### Problem Definition:

An *SBLS* of order  $n$  is an  $n \times n$  square grid in which:

- Each symbol appears exactly once in each row and column (*Latin square* structure; multiplication table of quasi-group).

*SBLS* of order 6



# Our Challenge Domain:

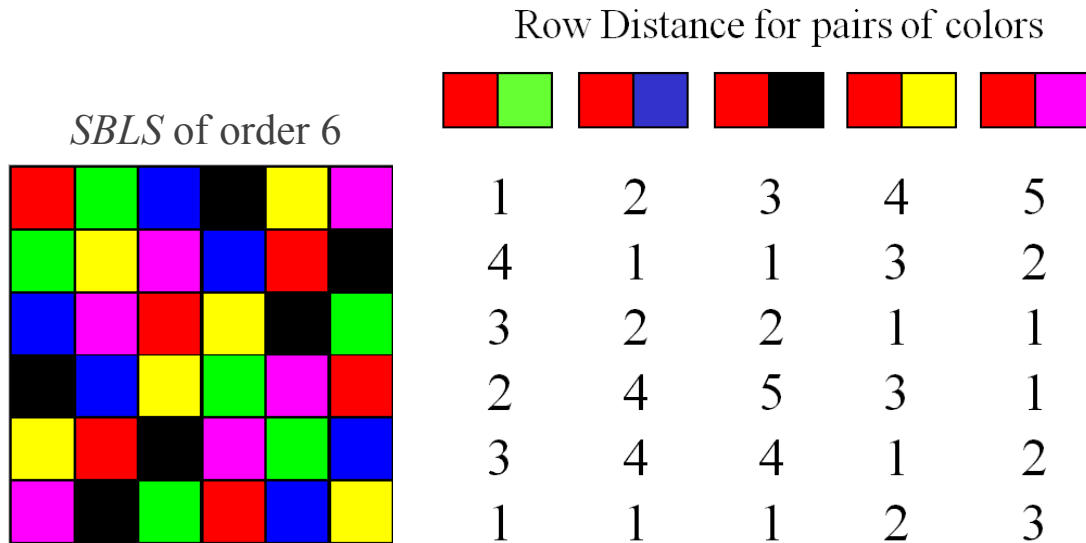
## The *Spatially-Balanced Latin Square (SBLs)* problem



### Problem Definition:

An *SBLs* of order  $n$  is an  $n \times n$  square grid in which:

- Each symbol appears exactly once in each row and column (*Latin square structure*; multiplication table of quasi-group).
- The average distance (column-wise) of a pair of symbols is the same for any pair (*Balanced structure*).



Total Row Distance (pair)

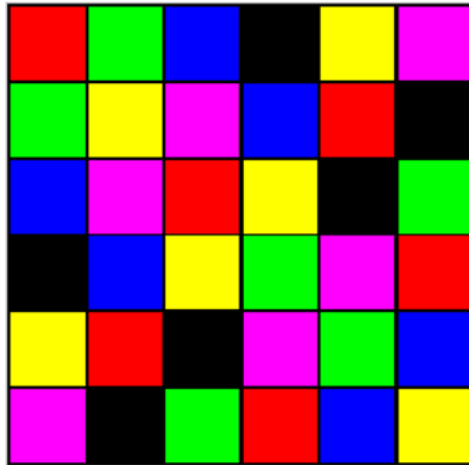
Average Row Distance (pair)



# The *Spatially-Balanced Latin Square (SBLs)* problem

## *SBLs* of order 6

Row Distance for pairs of colors



1	2	3	4	5
4	1	1	3	2
3	2	2	1	1
2	4	5	3	1
3	4	1	1	2
1	1	2	2	3

Total Row Distance (pair)	14	14	14	14	14
Average Row Distance (pair)	2.33	2.33	2.33	2.33	2.33



**14 (avg. 2.33) for all pairs. Use in “experimental design.”**

***In Computational Sustainability: to limit use of fertilizers.***

***Designs used by farmers in New York State!***

# Finding SBLS is very hard.

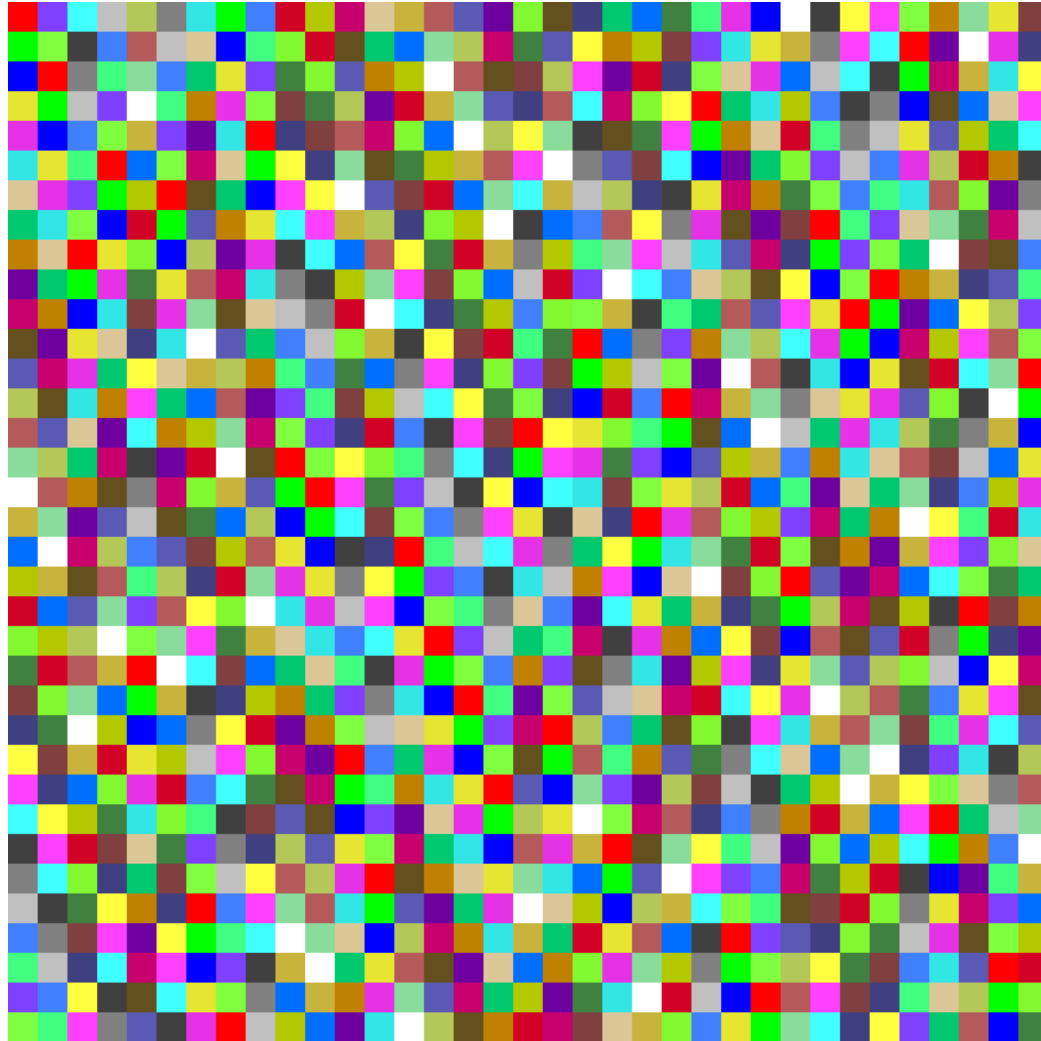
Approach	Order	Time (s)	Reference
Constraint Programming (CP)	9	241	[Gomes and Sellmann, CP'04]
IDWalk ( <i>metaheuristic</i> )	9	4.5	[Neveu et al., CP'04]
Self-symmetry-based Streamlined CP	14	5,434	[Gomes and Sellmann, CP'04]
Composition-based Streamlined CP	18	107K	[Gomes and Sellmann, CP'04]
Streamlined Local Search	<b>35</b>	<b>1.2M</b>	[Smith et al., IJCAI'05]

**Approach: use lots of problem specific structure to guide search.**

**Earlier known designs, only up to order 6.**

**The balancing condition (of each pair!) makes the problem \*very\* hard.**

# Largest ever found 35x35 --- 330 hrs cpu time.



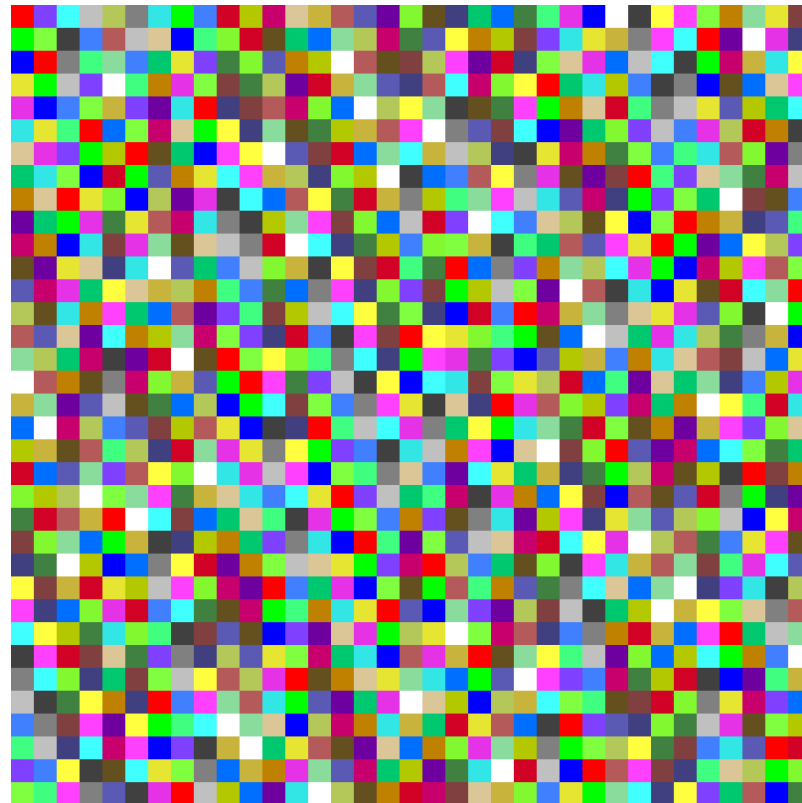


(here)

## Gedanken Experiment:

**Given that a constructive algorithm  
(given any size  $n$ ) exists, how would  
you find it? Hwk...**

**(Can use SAT solver.)**



# Outline

---

- Motivation
- Example Domain
- Proposed Framework
  - Overview of Streamlined Search
  - Taking advantage of Human Insights
  - Formal Description and Overview
  - Human-guided Streamlined Search
- Application to the *Spatially-balanced Latin square* problem
- Application to the *Weak Schur Number* problem
- Conclusions and Future work

# Outline

---

- Motivation
- Example Domain
- Proposed Framework
  - Overview of Streamlined Search
  - Taking advantage of Human Insights
  - Formal Description and Overview
  - GUI for Human-guided Streamlined Search
- Application to the *Spatially-balanced Latin square* problem
- Application to the *Weak Schur Number* problem
- Conclusions and Future work

# Proposed Framework: Overview of Streamlined Combinatorial Search

---



## Goal:

Exploit the **structure of some solutions** to **boost** the effectiveness of the **propagation mechanisms**.

## Underlying Observation:

When one insists on maintaining the **full solution set**, there is a **hard practical limit** on the effectiveness of **constraint propagation** methods. Often, there is **no compact representation** for all the solutions.

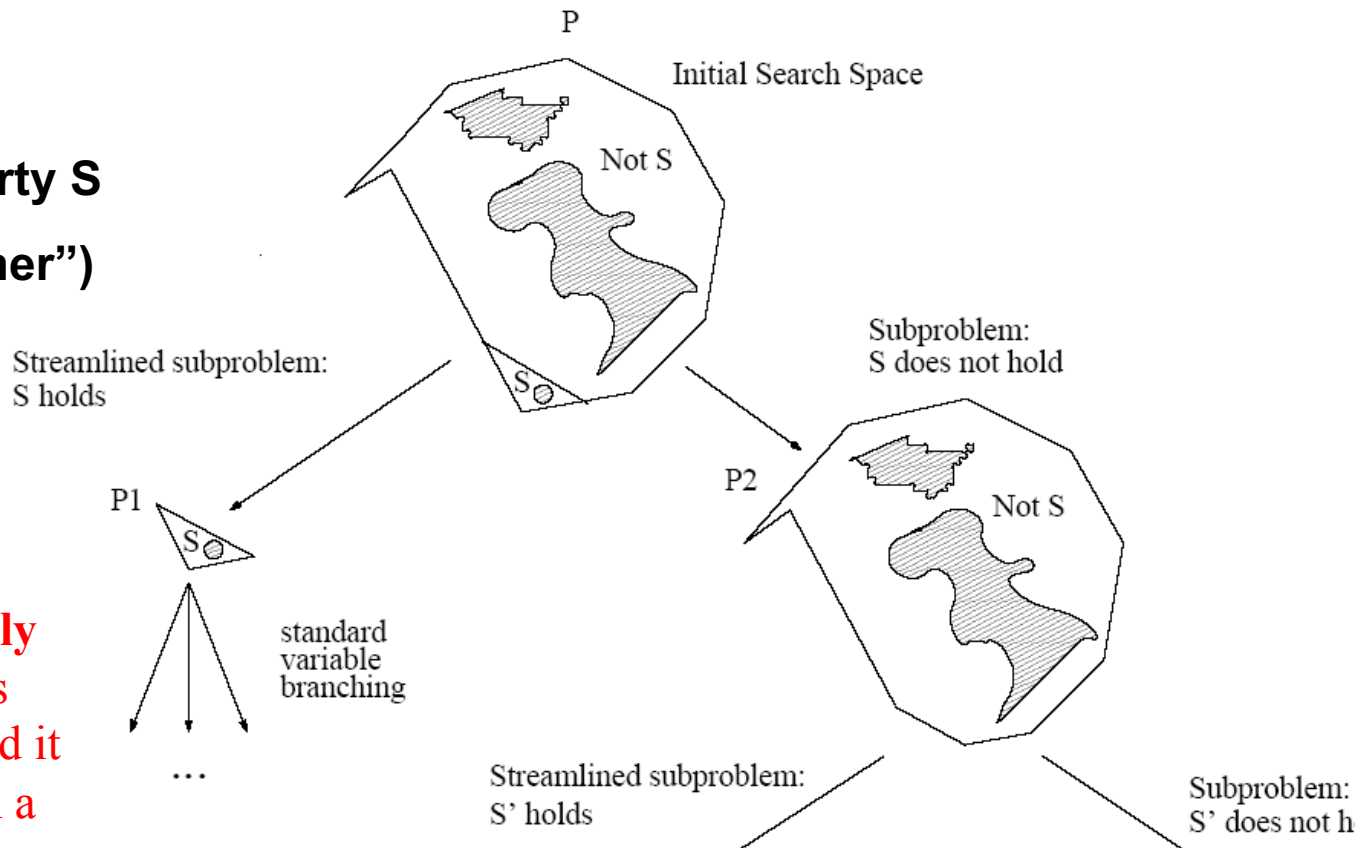
## Underlying Conjecture:

For many intricate **combinatorial problems** – if solutions exist – there will often be (highly) **regular ones**. **Can we find such solutions?**



# Overview of Streamlined Search

**Assert**  
**global property S**  
**(a “streamliner”)**

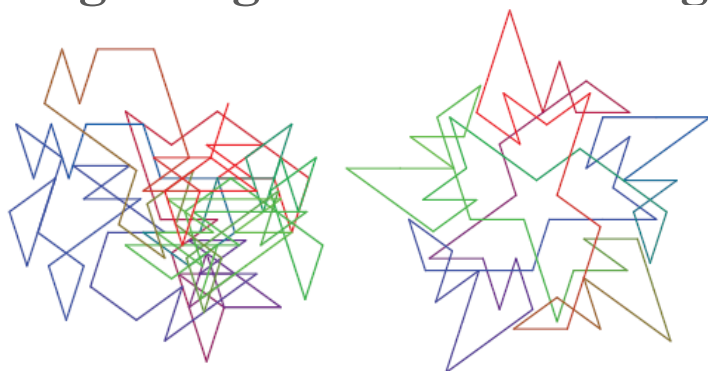


**$P_1$  is substantially smaller than its complement  $P_2$  and it will benefit from a stronger filtering thereafter.**

## Streamlined Search:

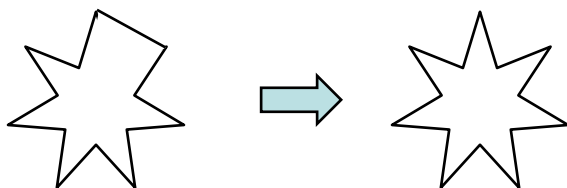
Strong **branching mechanisms** (by adding constraints based on **structure properties**) at **high levels** of the search tree.

## Recognizing Patterns and Regularities:



[Source: Marijn J.H. Heule, 2009, in work on van der Waerden numbers]

## Correcting Irregularities:



## Generalizing / Formalizing Regularities:

1	2	3
3	1	2
2	3	1

*Cyclic Latin square  
of order 3*



1	2	3	4
4	1	2	3
3	4	1	2
2	3	4	1

*Cyclic Latin square  
of order 4*

# Proposed Framework: Formal Description and Overview



```
 $\mathcal{O} \leftarrow \emptyset;$  // Conjectured streamliners  
 $\Gamma \leftarrow \emptyset;$  // Search streamliners  
 $\rho \leftarrow \rho_0;$  // Search parameter  
 $\mathcal{S} \leftarrow \emptyset;$  // Solutions found  
 $\tau \leftarrow false;$  // Timeout flag  
repeat  
   $Solve(P_\rho, \Gamma, t) \rightarrow (\mathcal{S}', \tau);$  // Search for new solutions  
  if  $\mathcal{S}' \cap \mathcal{S} \neq \emptyset$  then  
     $\mathcal{S} \leftarrow \mathcal{S} \cup \mathcal{S}';$  // Case 1: successful search  
     $Analyze(\mathcal{S}) \rightarrow \mathcal{O}';$  // Conjecture new streamliners  
     $\mathcal{O} \leftarrow \mathcal{O} \cup \mathcal{O}';$   
     $\rho \leftarrow \rho + 1;$   
  else if  $\tau$  is true then  
    Select  $\Gamma' \subseteq \mathcal{O};$  // Case 2: timed-out failed search  
     $\Gamma \leftarrow \Gamma \cup \Gamma';$  // Strengthen streamliners  
  else  
    Select  $\Gamma' \subseteq \Gamma;$  // Case 3: exhaustive failed search  
     $\Gamma \leftarrow \Gamma \setminus \Gamma';$  // Weaken streamliners  
     $\rho = max\{\rho : \mathcal{S}(\Gamma) \cap \mathcal{S}(P_\rho) \neq \emptyset\} + 1;$   
    Select  $\Gamma'' \subseteq \Gamma';$  // Find next parameter of interest  
     $\mathcal{O} \leftarrow \mathcal{O} \setminus \Gamma'';$  // Drop unpromising streamliners  
until  $\mathcal{O} = \emptyset;$ 
```

**Algorithm** : Discover-Construction procedure  
for a given problem  $P$ , with parameter set  $\rho$  and timeout  $t$ .

# Overview Approach --- Mimics Alg Discovery



```
 $\mathcal{O} \leftarrow \emptyset;$  // Conjectured streamliners  
 $\Gamma \leftarrow \emptyset;$  // Search streamliners  
 $\rho \leftarrow \rho_0;$  // Search parameter  
 $\mathcal{S} \leftarrow \emptyset;$  // Solutions found  
 $\tau \leftarrow \text{false};$  // Timeout flag  
repeat  
  2  $\text{Solve}(P_\rho, \Gamma, t) \rightarrow (\mathcal{S}', \tau);$  // Search for new solutions  
  if  $\mathcal{S}' \cap \mathcal{S} \neq \emptyset$  then  
    1  $\mathcal{S} \leftarrow \mathcal{S} \cup \mathcal{S}';$  // Case 1: successful search  
    1  $\text{Analyze}(\mathcal{S}) \rightarrow \mathcal{O}';$  // Conjecture new streamliners  
    4  $\mathcal{O} \leftarrow \mathcal{O} \cup \mathcal{O}';$   
    4  $\rho \leftarrow \rho + 1;$   
  else if  $\tau$  is true then  
    3  $\text{Select } \Gamma' \subseteq \mathcal{O};$  // Case 2: timed-out failed search  
     $\Gamma \leftarrow \Gamma \cup \Gamma';$  // Strengthen streamliners  
  else  
     $\text{Select } \Gamma' \subseteq \Gamma;$  // Case 3: exhaustive failed search  
     $\Gamma \leftarrow \Gamma \setminus \Gamma';$  // Weaken streamliners  
     $\rho = \max\{\rho : \mathcal{S}(\Gamma) \cap \mathcal{S}(P_\rho) \neq \emptyset\} + 1;$   
     $\text{Select } \Gamma'' \subseteq \Gamma';$  // Find next parameter of interest  
     $\mathcal{O} \leftarrow \mathcal{O} \setminus \Gamma'';$  // Drop unpromising streamliners  
until  $\mathcal{O} = \emptyset;$ 
```

**Algorithm** : Discover-Construction procedure  
for a given problem  $P$ , with parameter set  $\rho$  and timeout  $t$ .

- 1 Analyze smaller size solutions, and **conjecture potential regularities** in the solutions. (Human insight.)
- 2 Validate through **streamlining** the observed regularities.
- 3 If the streamlined search **does not give a larger size solution**, the proposed regularity is quite likely **accidental** and one looks for a new pattern in the small scale solutions.
- 4 Otherwise, one proceeds by generating a number of **new solutions** that all contain the proposed **structural regularity** and are used to expand the solution set and to **reveal new regularities**.

# Proposed Framework for Human-guided Streamlined Search



## Overview on the *SBSL* problem

---

Search Parameters

$$n=3$$

$$\Gamma=\{\}$$

Conjectured Streamliners

Start with first order of interest  
( $n=3$ ) and no streamliners ( $\Gamma=\{\}$ )

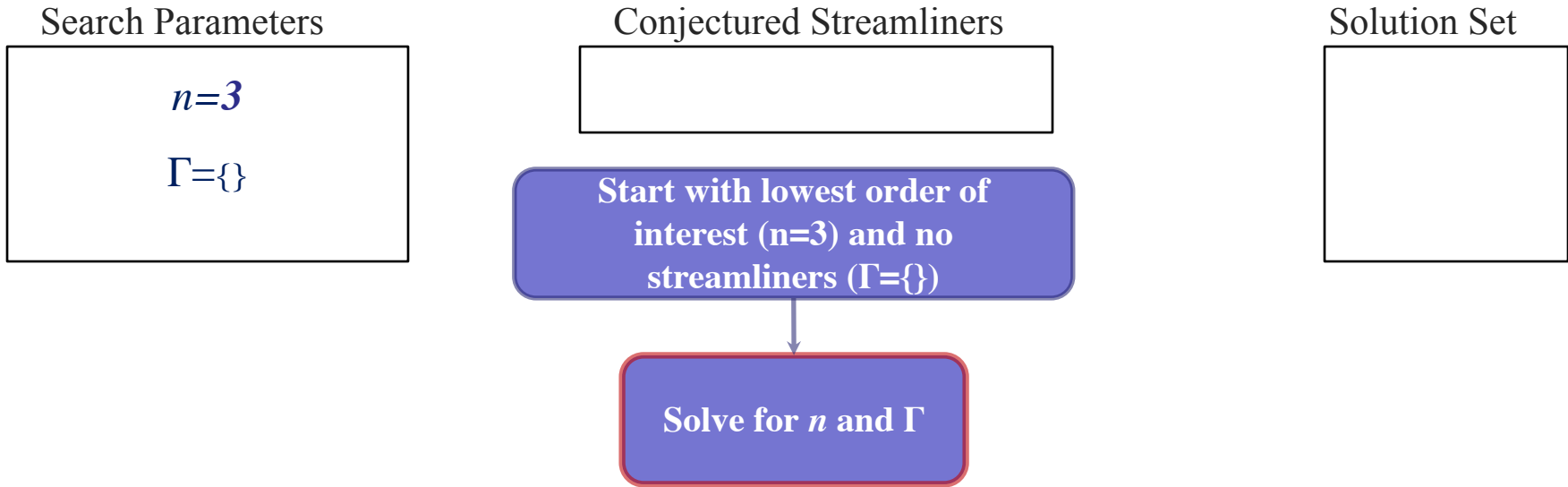
Solution Set

# Proposed Framework for Human-guided Streamlined Search



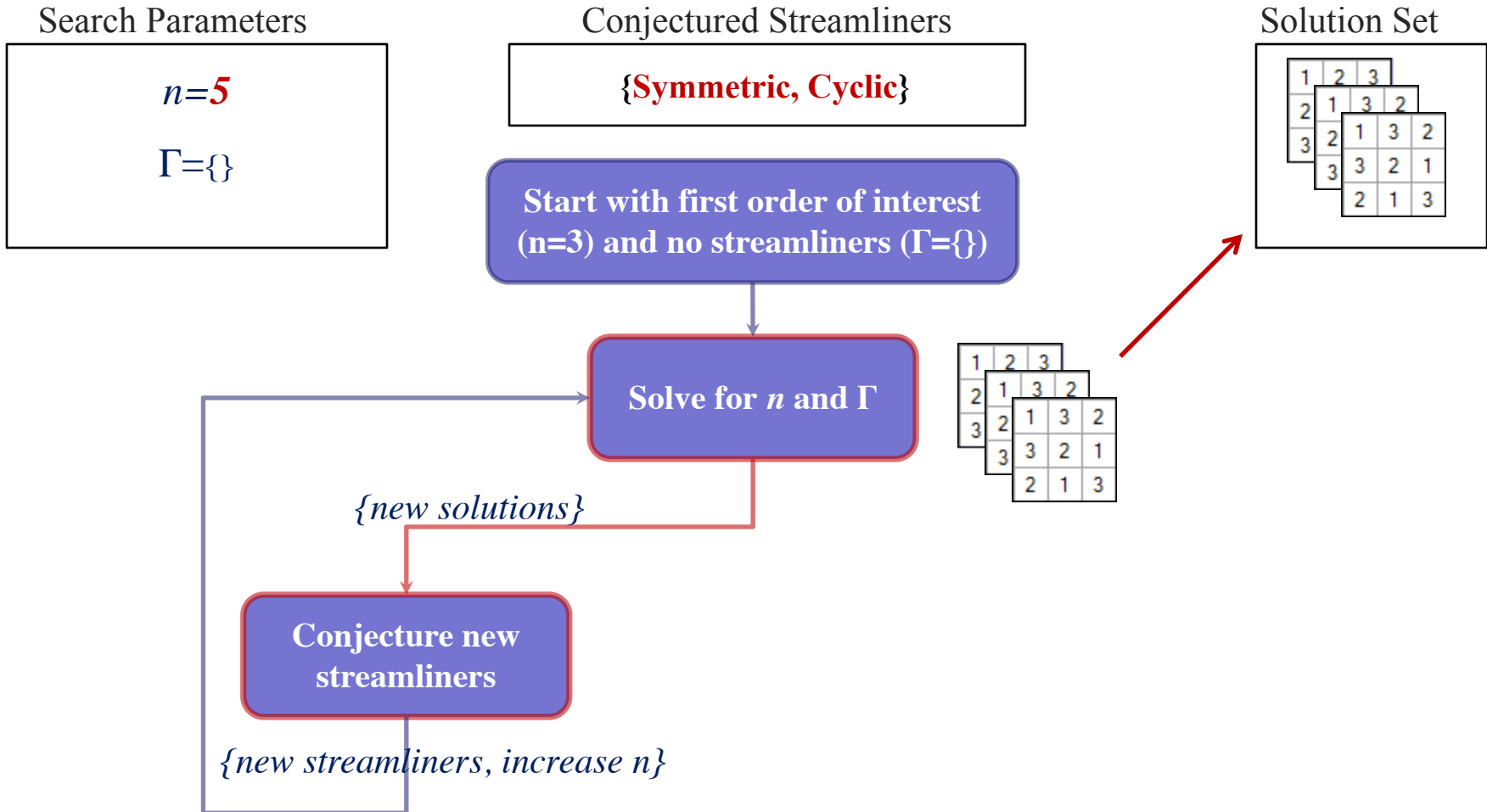
## Overview on the *SBS* problem

---



# Proposed Framework for Human-guided Streamlined Search

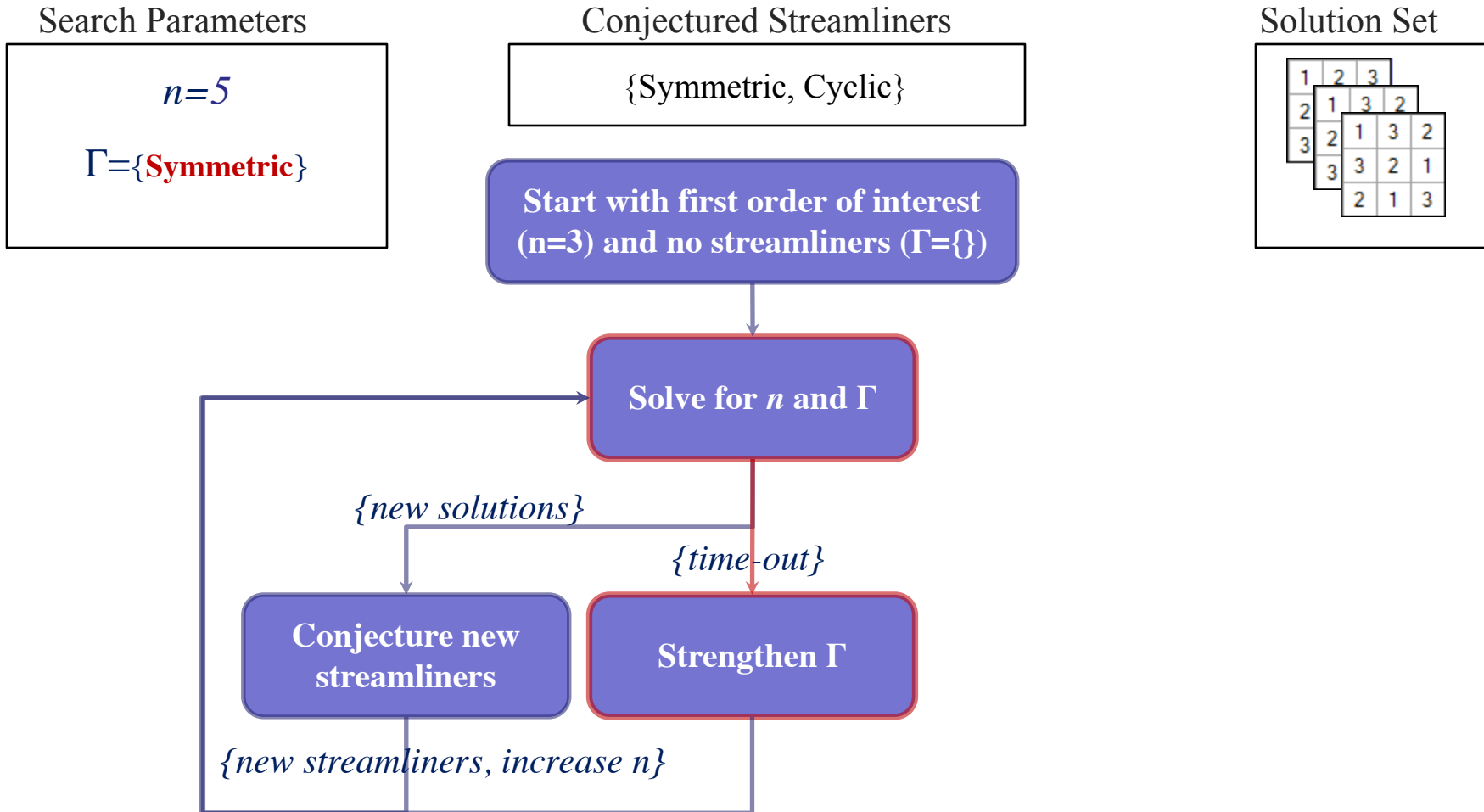
## Overview on the *SBLS* problem



# Proposed Framework for Human-guided Streamlined Search



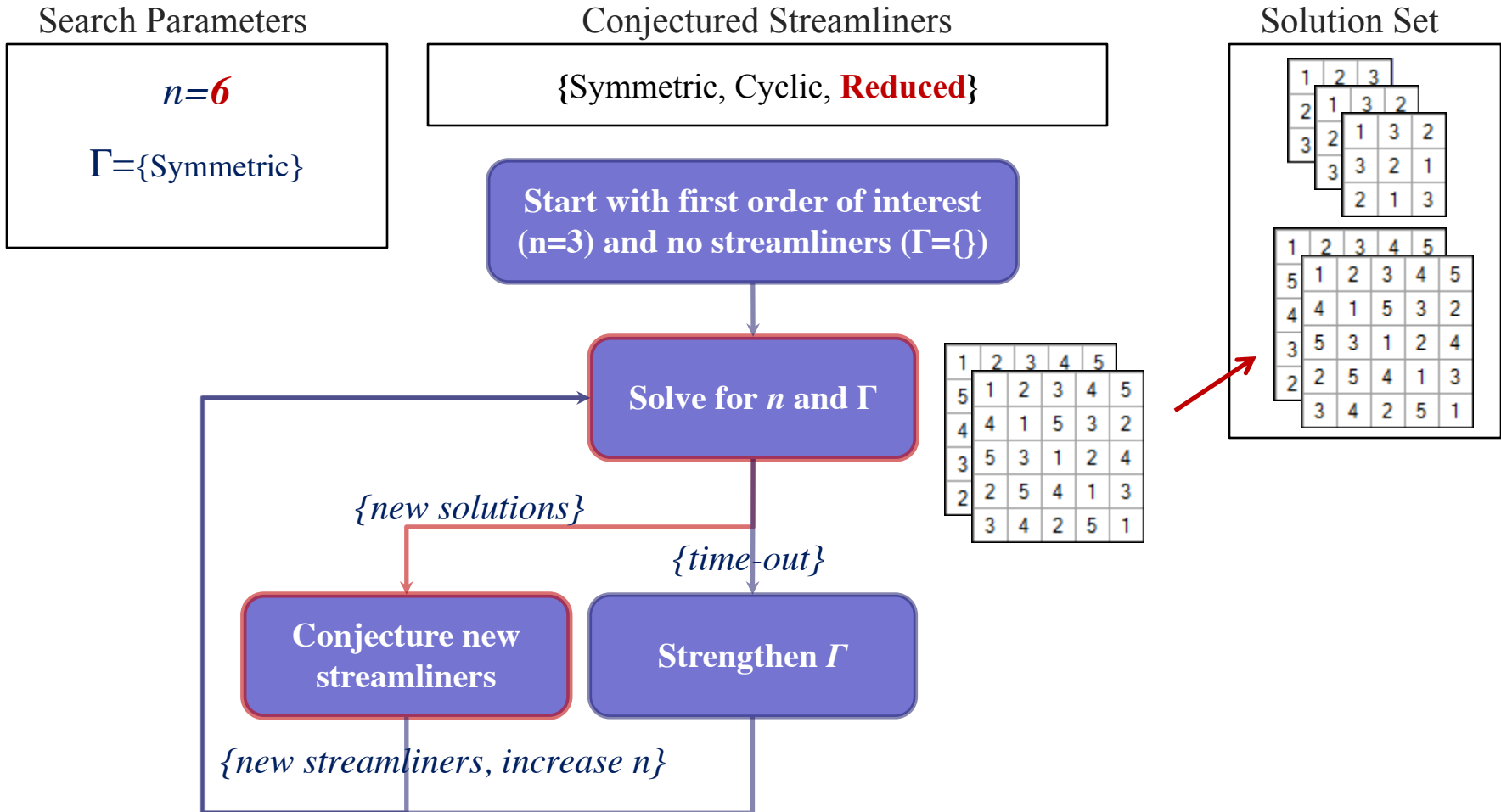
## Overview on the *SBL* problem





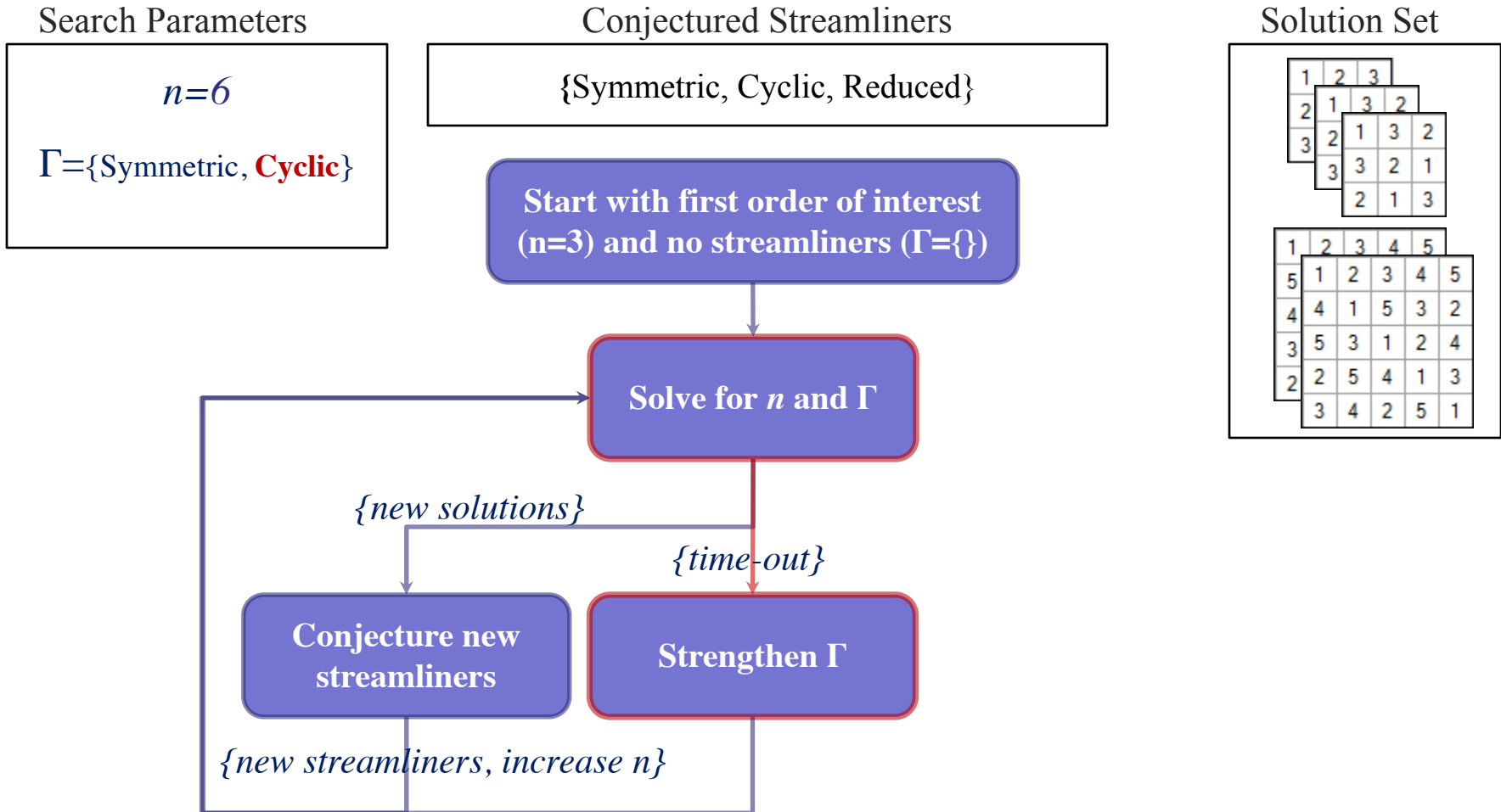
# Proposed Framework for Human-guided Streamlined Search

## Overview on the *SBL* problem



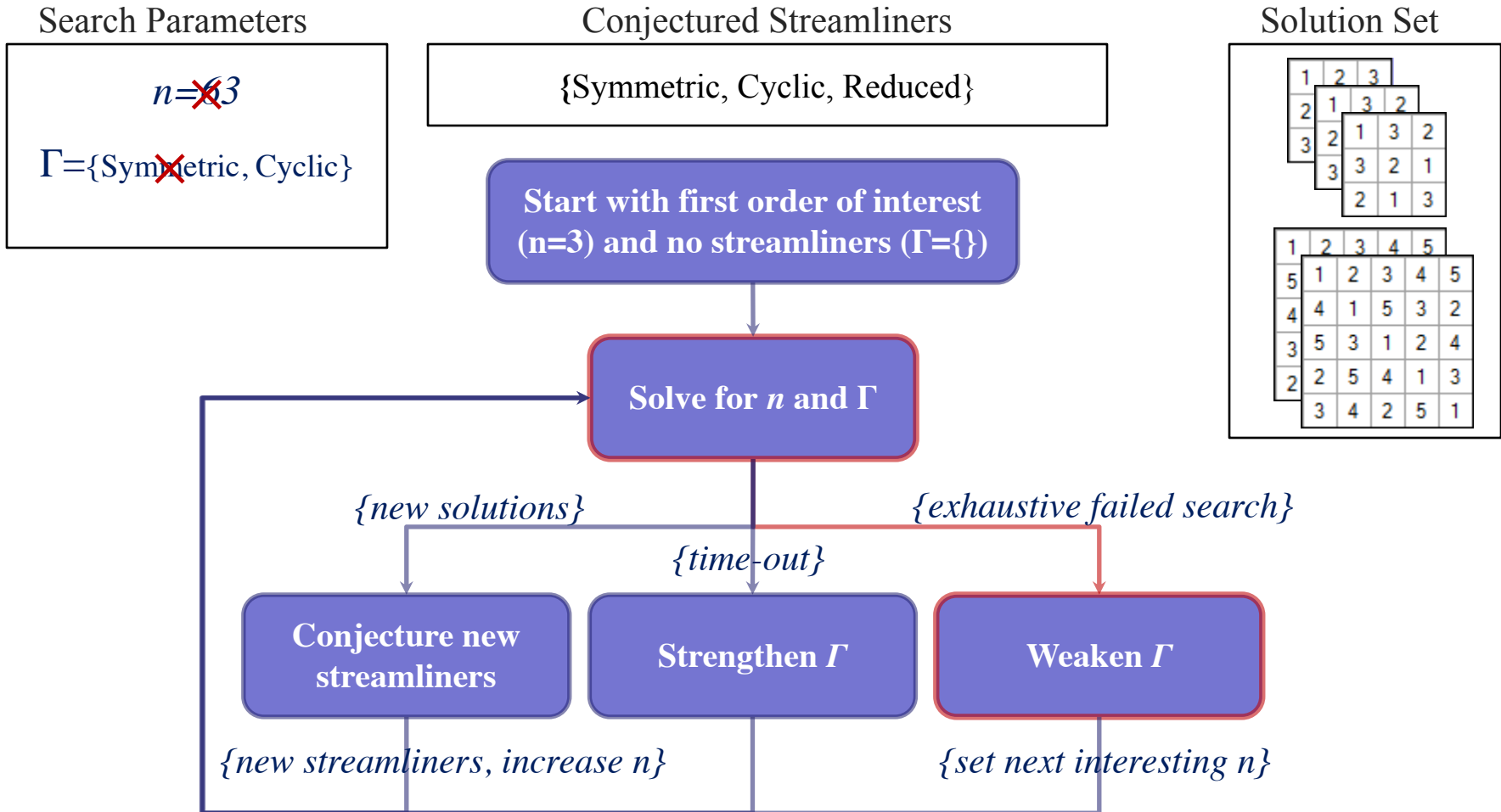
# Proposed Framework for Human-guided Streamlined Search

## Overview on the *SBL* problem



# Proposed Framework for Human-guided Streamlined Search

## Overview on the *SBL* problem



# GUI for Human-guided Streamlined Search

Constructive Procedures Discovery Tool

File Edit Help

### Solutions found

Streamliner \ Parameter	3	5	6	8	9	11	12
Any	12	5760	8736	238	411	9	6
Reduced	1	2	14	12	1	1	0
Symmetry	6	240	8640	12	1	1	0
Columns 2 and n	1	6	1	2	1	1	0
Cyclic	6	40	96	226	410	8	6

### Selected Solutions

5:1464
5:3194
5:4526
5:4983
5:5502
5:5572
6:10
8:10
8:11

1	2	3	4	5
2	4	5	3	1
3	5	2	1	4
4	3	1	5	2
5	1	4	2	3

1	2	3	4	5	6
2	4	6	5	3	1
3	6	4	1	2	5
4	5	1	3	6	2
5	3	2	6	1	4
6	1	5	2	4	3

1	2	3	4	5	6	7	8
2	4	6	8	7	5	3	1
3	6	8	5	2	1	4	7
4	8	5	1	3	7	6	2
5	7	2	3	8	4	1	6
6	5	1	7	4	2	8	3
7	3	4	6	1	8	2	5
8	1	7	2	6	3	5	4

### Streamliner Editor

Streamliner Name: Symmetry

Evaluation Function

```
for(int i=0; i<n; i++){
  for(int j=0; j<n; j++){
    if( a[i][j] != a[j][i] ){
      return false;
    }
  }
}
return true;
```

Constraint Post Function

```
for(int i=0; i<n; i++){
  for(int j=0; j<n; j++){
    cp.Add(cp.Eq(a[i][j],a[j][i]));
  }
}
```

Cancel OK

### 1 - Select Streamliner Combination

Reduced  
Symmetry  
Columns 2 and n  
Cyclic

Reduced  
Symmetry  
Columns 2 and n

### 2 - Set Parameters

Parameter n: 11  
Parameter k:   
Create New Streamliner:   
Name:

### 3 - Perform search

Time Limit: 60  
Click to Run:

### Search Stats

New solutions found: 1  
Total # of solutions for this order: 9  
Total # of solutions: 15172

# GUI for Human-guided Streamlined Search

Constructive Procedures Discovery Tool

File Edit Help

### Solutions found

Streamliner \ Parameter	3	5	6	8	9	11	12
Any	12	5760	8736	238	411	9	6
Reduced	1	2	14	12	1	1	0
Symmetry	6	240	8640	12	1	1	0
Columns 2 and n	1	6	1	2	1	1	0
Cyclic	6	40	96	226	410	8	6

### Selected Solutions

5:1464
5:3194
5:4526
5:4983
5:5502
5:5572
6:10
8:10
8:11

1	2	3	4	5
2	4	5	3	1
3	5	2	1	4
4	3	1	5	2
5	1	4	2	3

1	2	3	4	5	6
2	4	6	5	3	1
3	6	4	1	2	5
4	5	1	3	6	2
5	3	2	6	1	4
6	1	5	2	4	3

1	2	3	4	5	6	7	8
2	4	6	8	7	5	3	1
3	6	8	5	2	1	4	7
4	8	5	1	3	7	6	2
5	7	2	3	8	4	1	6
6	5	1	7	4	2	8	3
7	3	4	6	1	8	2	5
8	1	7	2	6	3	5	4

### Streamliner Editor

Streamliner Name: Symmetry

Evaluation Function

```
for(int i=0; i<n; i++){
  for(int j=0; j<n; j++){
    if( a[i][j] != a[j][i] ){
      return false;
    }
  }
}
return true;
```

Constraint Post Function

```
for(int i=0; i<n; i++){
  for(int j=0; j<n; j++){
    cp.Add(cp.Eq(a[i][j],a[j][i]));
  }
}
```

Cancel OK

### 1 - Select Streamliner Combination

Reduced  
Symmetry  
Columns 2 and n  
Cyclic

Reduced  
Symmetry  
Columns 2 and n

### 2 - Set Parameters

Parameter n: 11  
Parameter k:   
Create New Streamliner:   
Name:

### 3 - Perform Search

Time Limit: 60  
Click to Run: Streamline

### Search Stats

New solutions found: 9  
Total # of solutions for this order: 9  
Total # of solutions: 15172

Conjecture new streamliners

# GUI for Human-guided Streamlined Search

Constructive Procedures Discovery Tool

File Edit Help

### Solutions found

Streamliner \ Parameter	3	5	6	8	9	11	12
Any	12	5760	8736	238	411	9	6
Reduced	1	2	14	12	1	1	0
Symmetry	6	240	8640	12	1	1	0
Columns 2 and n	1	6	1	2	1	1	0
Cyclic	6	40	96	226	410	8	6

### Selected Solutions

5:1464	1 2 3 4 5	1 2 3 4 5 6	1 2 3 4 5 6 7 8
5:3194	2 4 5 3 1	2 4 6 5 3 1	2 4 6 8 7 5 3 1
5:4526	3 5 2 1 4	3 6 4 1 2 5	3 6 8 5 2 1 4 7
5:4983	4 3 1 5 2	4 5 1 3 6 2	4 8 5 1 3 7 6 2
5:5502	5 1 4 2 3	5 3 2 6 1 4	5 7 2 3 8 4 1 6
5:5572		6 1 5 2 4 3	6 5 1 7 4 2 8 3
6:10			7 3 4 5 6 2 1 8
6:10			8 1 7 2 6 3 5 4

Select streamliners, parameters, and perform the search

#### 1 - Select Streamliner Combination

Reduced  
 Symmetry  
 Columns 2 and n  
 Cyclic

#### 2 - Set Parameters

Parameter n:   
 Parameter k:   
 Create New Streamliner  
 Name:

#### 3 - Perform search

Time Limit:   
 Click to Run:

#### Search Stats

New solutions found:   
 Total # of solutions for this order:   
 Total # of solutions:

#### Streamliner Editor

Streamliner Name:

Evaluation Function

```

for(int i=0; i<n; i++){
  for(int j=0; j<n; j++){
    if( a[i][j] != a[j][i] ){
      return false;
    }
  }
}
return true;
        
```

Constraint Post Function

```

for(int i=0; i<n; i++){
  for(int j=0; j<n; j++){
    cp.Add(cp.Eq(a[i][j],a[j][i]));
  }
}
        
```

# GUI for Human-guided Streamlined Search

Constructive Procedures Discovery Tool

File Edit **Select solutions**

Streamliner \ Parameter	3	5	6	8	9	11	12
Any	12	5760	8736	238	411	9	6
Reduced	1	2	14	12	1	1	0
Symmetry	6	240	8640	12	1	1	0
Columns 2 and n	1	6	1	2	1	1	0
Cyclic	6	40	96	226	410	8	6

**Streamliner Editor**

Streamliner Name: Symmetry

Evaluation Function

```
for(int i=0; i<n; i++){
  for(int j=0; j<n; j++){
    if( a[i][j] != a[j][i] ){
      return false;
    }
  }
}
return true;
```

Constraint Post Function

```
for(int i=0; i<n; i++){
  for(int j=0; j<n; j++){
    cp.Add(cp.Eq(a[i][j],a[j][i]));
  }
}
```

Selected Solutions

- 5:1464
- 5:3194
- 5:4526
- 5:4983
- 5:5502
- 5:5572
- 6:10
- 8:10
- 8:11

1	2	3	4	5
2	4	5	3	1
3	5	2	1	4
4	3	1	5	2
5	1	4	2	3

1	2	3	4	5	6
2	4	6	5	3	1
3	6	4	1	2	5
4	5	1	3	6	2
5	3	2	6	1	4
6	1	5	2	4	3

1	2	3	4	5	6	7	8
2	4	6	8	7	5	3	1
3	6	8	5	2	1	4	7
4	8	5	1	3	7	6	2
5	7	2	3	8	4	1	6
6	5	1	7	4	2	8	3
7	3	4	6	1	8	2	5
8	1	7	2	6	3	5	4

1 - Select Streamliner Combination

Reduced  
Symmetry  
Columns 2 and n  
Cyclic

Reduced  
Symmetry  
Columns 2 and n

2 - Set Parameters

Parameter n: 11  
Parameter k:   
Create New Streamliner:   
Name:

3 - Perform search

Time Limit: 60  
Click to Run:

Search Stats

New solutions found: 1  
Total # of solutions for this order: 9  
Total # of solutions: 15172

# GUI for Human-guided Streamlined Search

Constructive Procedures Discovery Tool

File Edit Help

### Solutions found

Streamliner \ Parameter	3	5	6	8	9	11	12
Any	12	5760	8736	238	411	9	6
Reduced	1	2	14	12	1	1	0
Symmetry	6	240	8640	12	1	1	0
Columns 2 and n	1	6	1	2	1	1	0
Cyclic	96			226	410	8	6

**Analyze solutions**

### Selected Solutions

- 5:1464
- 5:3194
- 5:4526
- 5:4983
- 5:5502
- 5:5572
- 6:10
- 8:10
- 8:11

1	2	3	4	5
2	4	5	3	1
3	5	2	1	4
4	3	1	5	2
5	1	4	2	3

1	2	3	4	5	6
2	4	6	5	3	1
3	6	4	1	2	5
4	5	1	3	6	2
5	3	2	6	1	4
6	1	5	2	4	3

1	2	3	4	5	6	7	8
2	4	6	8	7	5	3	1
3	6	8	5	2	1	4	7
4	8	5	1	3	7	6	2
5	7	2	3	8	4	1	6
6	5	1	7	4	2	8	3
7	3	4	6	1	8	2	5
8	1	7	2	6	3	5	4

### Streamliner Editor

Streamliner Name: Symmetry

Evaluation Function

```
for(int i=0; i<n; i++){
  for(int j=0; j<n; j++){
    if( a[i][j] != a[j][i] ){
      return false;
    }
  }
}
return true;
```

Constraint Post Function

```
for(int i=0; i<n; i++){
  for(int j=0; j<n; j++){
    cp.Add(cp.Eq(a[i][j],a[j][i]));
  }
}
```

Cancel OK

### 1 - Select Streamliner Combination

Reduced  
Symmetry  
Columns 2 and n  
Cyclic

Reduced  
Symmetry  
Columns 2 and n

### 2 - Set Parameters

Parameter n: 11  
Parameter k:   
Create New Streamliner:   
Name:

### 3 - Perform search

Time Limit: 60  
Click to Run:

### Search Stats

New solutions found: 1  
Total # of solutions for this order: 9  
Total # of solutions: 15172



# Outline

---

- Motivation
- Example Domain
- Proposed Framework
- Application to the *Spatially-Balanced Latin square* problem
  - Successful Streamliners
  - Constructive Procedure 1
  - Constructive Procedure 2
- Application to the *Weak Schur Number* problem
- Conclusions and Future work

# Application to the *SBLs* problem: Construction 1

## *Successful Key Streamliners:*

{Diagonal symmetry, Reduced form, Assignments of columns 2 and  $n$ , Multiples of  $i$  in row  $i$ , Second sequence decreasing}

Streamliners	5	6	8	9	11	14
$\Gamma_1 = \emptyset$	<b>5760</b>	15878	-	-	-	-
$\Gamma_2 = \Gamma_1 \cup \{\text{Symmetric}\}$	<b>240</b>	8447	714	43	-	-
$\Gamma_3 = \Gamma_2 \cup \{\text{Reduced}\}$	<b>2</b>	<b>14</b>	<b>14</b>	51	-	-
$\Gamma_4 = \Gamma_3 \cup \{\text{Columns 2 \& } n\}$	<b>1</b>	<b>1</b>	<b>2</b>	<b>1</b>	<b>1</b>	-
$\Gamma_5 = \Gamma_4 \cup \{\text{Multiples of } i\}$	<b>1</b>	<b>1</b>	<b>2</b>	<b>1</b>	<b>1</b>	<b>1</b>

*Fig: Number of SBLs generated in 60 seconds, by order and streamliners (Bold indicates exhaustive search).*

# Application to the *SBLs* problem: Construction 1

1	2	3	4	5	6	7	8
2	4	6	8	7	5	3	1
3	6	8	5	2	1	4	7
4	8	5	1	3	7	6	2
5	7	2	3	8	4	1	6
6	5	1	7	4	2	8	3
7	3	4	6	1	8	2	5
8	1	7	2	6	3	5	4

1	2	3	4	5	6	7	8	9
2	4	6	8	9	7	5	3	1
3	6	9	7	4	1	2	5	8
4	8	7	3	1	5	9	6	2
5	9	4	1	6	8	3	2	7
6	7	1	5	8	2	4	9	3
7	5	2	9	3	4	8	1	6
8	3	5	6	2	9	1	7	4
9	1	8	2	7	3	6	4	5

1	2	3	4	5	6	7	8	9	10	11
2	4	6	8	10	11	9	7	5	3	1
3	6	9	11	8	5	2	1	4	7	10
4	8	11	7	3	1	5	9	10	6	2
5	10	8	3	2	7	11	6	1	4	9
6	11	5	1	7	10	4	2	8	9	3
7	9	2	5	11	4	3	10	6	1	8
8	7	1	9	6	2	10	5	3	11	4
9	5	4	10	1	8	6	3	11	2	7
10	3	7	6	4	9	1	11	2	8	5
11	1	10	2	9	3	8	4	7	5	6

1	2	3	4	5	6	7	8	9	10	11	12	13	14
2	4	6	8	10	12	14	13	11	9	7	5	3	1
3	6	9	12	14	11	8	5	2	1	4	7	10	13
4	8	12	13	9	5	1	3	7	11	14	10	6	2
5	10	14	9	4	1	6	11	13	8	3	2	7	12
6	12	11	5	1	7	13	10	4	2	8	14	9	3
7	14	8	1	6	13	9	2	5	12	10	3	4	11
8	13	5	3	11	10	2	6	14	7	1	9	12	4
9	11	2	7	13	4	5	14	6	3	12	8	1	10
10	9	1	11	8	2	12	7	3	13	6	4	14	5
11	7	4	14	3	8	10	1	12	6	5	13	2	9
12	5	7	10	2	14	3	9	8	4	13	1	11	6
13	3	10	6	7	9	4	12	1	14	2	11	5	8
14	1	13	2	12	3	11	4	10	5	9	6	8	7

# Application to the *SBL*s problem: Construction 1

1	2	3	4	5	6	7	8	9	10	11	12	13	14
2	4	6	8	10	12	14	13	11	9	7	5	3	1
3	6	9	12	14	11	8	5	2	1	4	7	10	13
4	8	12	13	9	5	1	3	7	11	14	10	6	2
5	10	14	9	4	1	6	11	13	8	3	2	7	12
6	12	11	5	1	7	13	10	4	2	8	14	9	3
7	14	8	1	6	13	9	2	5	12	10	3	4	11
8	13	5	3	11	10	2	6	14	7	1	9	12	4
9	11	2	7	13	4	5	14	6	3	12	8	1	10
10	9	1	11	8	2	12	7	3	13	6	4	14	5
11	7	4	14	3	8	10	1	12	6	5	13	2	9
12	5	7	10	2	14	3	9	8	4	13	1	11	6
13	3	10	6	7	9	4	12	1	14	2	11	5	8
14	1	13	2	12	3	11	4	10	5	9	6	8	7

a	a+i	a	a-i
---	-----	---	-----

**Approach reveals patterns. Final step is figuring out boundary conditions.**

# Application to the *SBLs* problem: Construction 1

1	2	3	4	5	6	7	8	9	10	11	12	13	14
2	4	6	8	10	12	14	13	11	9	7	5	3	1
3	6	9	12	14	11	8	5	2	1	4	7	10	13
4	8	12	13	9	5	1	3	7	11	14	10	6	2
5	10	14	9	4	1	6	11	13	8	3	2	7	12
6	12	11	5	1	7	13	10	4	2	8	14	9	3
7	14	8	1	6	13	9	2	5	12	10	3	4	11
8	13	5	3	11	10	2	6	14	7	1	9	12	4
9	11	2	7	13	4	5	14	6	3	12	8	1	10
10	9	1	11	8	2	12	7	3	13	6	4	14	5
11	7	4	14	3	8	10	1	12	6	5	13	2	9
12	5	7	10	2	14	3	9	8	4	13	1	11	6
13	3	10	6	7	9	4	12	1	14	2	11	5	8
14	1	13	2	12	3	11	4	10	5	9	6	8	7

a	a+i	a	a-i
---	-----	---	-----

```

for row  $i = 1, \dots, N$  do
     $k = 1;$  // Sequence number
     $j = 1;$  // Column index
     $a_{i,j} = i;$  // First symbol of row  $i$ 
    while  $j < N$  do
        if  $k$  is odd then // Odd sequence
            while  $a_{i,j} + i \leq N$  and  $j < N$  do
                 $a_{i,j+1} = a_{i,j} + i;$ 
                 $j = j + 1;$ 
            else // Even sequence
                while  $a_{i,j} - i \geq 1$  and  $j < N$  do
                     $a_{i,j+1} = a_{i,j} - i;$ 
                     $j = j + 1;$ 
            if  $j < N$  then // Switch sequence
                if  $k$  is odd then
                     $a_{i,j+1} = 2N + 1 - i - a_{i,j};$ 
                else
                     $a_{i,j+1} = i - a_{i,j};$ 
                 $k = k + 1;$ 
                 $j = j + 1;$ 

```

**Algorithm :** SBLs-sequence procedure for SBLs of order  $N$ , when  $2N + 1$  is prime.

# Application to the *SBLS* problem: Construction 1

1	2	3	4	5	6	7	8	9	10	11	12	13	14
2	4	6	8	10	12	14	13	11	9	7	5	3	1
3	6	9	12	14	11	8	5	2	1	4	7	10	13
4	8	12	13	9	5	1	3	7	11	14	10	6	2
5	10	14	9	4	1	6	11	13	8	3	2	7	12
6	12	11	5	1	7	13	10	4	2	8	14	9	3
7	14	8	1	6	13	9	2	5	12	10	3	4	11
8	13	5	3	11	10	2	6	14	7	1	9	12	4
9	11	2	7	13	4	5	14	6	3	12	8	1	10
10	9	1	11	8	2	12	7	3	13	6	4	14	5
11	7	4	14	3	8	10	1	12	6	5	13	2	9
12	5	7	10	2	14	3	9	8	4	13	1	11	6
13	3	10	6	7	9	4	12	1	14	2	11	5	8
14	1	13	2	12	3	11	4	10	5	9	6	8	7

a	a+i	a	a-i
---	-----	---	-----

```

for row  $i = 1, \dots, N$  do
     $k = 1;$  // Sequence number
     $j = 1;$  // Column index
     $a_{i,j} = i;$  // First symbol of row  $i$ 
    while  $j < N$  do
        if  $k$  is odd then // Odd sequence
            while  $a_{i,j} + i \leq N$  and  $j < N$  do
                 $a_{i,j+1} = a_{i,j} + i;$ 
                 $j = j + 1;$ 
            else // Even sequence
                while  $a_{i,j} - i \geq 1$  and  $j < N$  do
                     $a_{i,j+1} = a_{i,j} - i;$ 
                     $j = j + 1;$ 
            if  $j < N$  then // Switch sequence
                if  $k$  is odd then
                     $a_{i,j+1} = 2N + 1 - i - a_{i,j};$ 
                else
                     $a_{i,j+1} = i - a_{i,j};$ 
                 $k = k + 1;$ 
                 $j = j + 1;$ 

```

**Algorithm :** SBLS-sequence procedure for SBLS of order  $N$ , when  $2N + 1$  is prime.

Proof of Correctness in [R. Le Bras, et al., *Polynomial Time Construction for Spatially Balanced Latin Squares*, 2012]

# Application to the *SBLs* problem: Construction 2

1	2	6	3	9	7	13	4	11	10	12	8	5	14
2	3	7	4	10	8	14	5	12	11	13	9	6	1
3	4	8	5	11	9	1	6	13	12	14	10	7	2
4	5	9	6	12	10	2	7	14	13	1	11	8	3
5	6	10	7	13	11	3	8	1	14	2	12	9	4
6	7	11	8	14	12	4	9	2	1	3	13	10	5
7	8	12	9	1	13	5	10	3	2	4	14	11	6
8	9	13	10	2	14	6	11	4	3	5	1	12	7
9	10	14	11	3	1	7	12	5	4	6	2	13	8
10	11	1	12	4	2	8	13	6	5	7	3	14	9
11	12	2	13	5	3	9	14	7	6	8	4	1	10
12	13	3	14	6	4	10	1	8	7	9	5	2	11
13	14	4	1	7	5	11	2	9	8	10	6	3	12
14	1	5	2	8	6	12	3	10	9	11	7	4	13

```

 $c_{1,1} = 1;$  // Generate 1st row of the conjugate
for column  $j = 2, \dots, N$  do // Observed pattern 1, 2, 4, ...
    if  $2c_{1,j-1} \leq N$  then
         $c_{1,j} = 2c_{1,j-1};$ 
    else
         $c_{1,j} = 2N + 1 - 2c_{1,j-1};$ 
for row  $i = 2, \dots, N$  do // Subsequent rows
     $c_{i,1} = c_{i-1,N};$  // Shifted version of previous
    for column  $j = 2, \dots, N$  do
         $c_{i,j} = c_{i-1,j-1};$ 
for row  $i = 1, \dots, N$  do // Generate SBLs from conjugate
    for column  $j = 1, \dots, N$  do
         $a_{i,c_{i,j}} = j;$ 
    
```

**Algorithm:** SBLs-Cyclic procedure.

## **Incremental approach in terms of size is necessary:**

**Need to reveal structure by generating increasingly larger solutions using increased streamlining (and thus, structure).**

*True solution structure only becomes visible around order 14.*

**At that size**

- (1) Can't generate any solution of that size without strong streamlining,**
- (2) But, even if we could, arbitrary solutions will not show structure.**



# Observation II

---



**We conjecture that many hard combinatorial design problems have effective constructive procedures. (E.g., quasigroup existence problems, code design, Ramsey graphs (numbers), van der Waerden, Schur numbers etc.).**

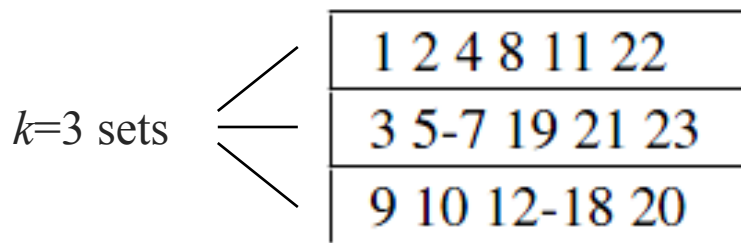
**So far, such problems have been tackled with combinatorial search (SAT solvers, equational thm. provers) on a per instance basis. Does resolve open questions in combinatorial design (e.g. new lower-bounds) but does not reveal the general solution structure.**

**E.g. Slaney et al. 1993, Herwig et al. 2007, Heule and Walsh 2010  
Eliahou et al. 2012, Fujita et al. 2013.**

# Example: the *Weak Schur* problem

## Problem Definition:

- A set is (*weakly*) **sum free** if for any two (*distinct*) elements of this set, their sum does not belong to the set.
- The *Weak Schür Number* of order  $k$ ,  $WS(k)$ , is the largest integer  $n$  for which there exists a partition of  $[1,n]$  into  $k$  weakly sum-free sets.



*Each of the 3 sets is such that, for any 2 elements of the set, their sum does not belong to the same set.*

**Fig: Partition of  $[1,23]$  into 3 weakly sum-free sets, proving  $WS(3) \geq 23$**

# Application to the *Weak Schur* problem

## Best known lower bounds:

Approach	$WS(5)$	$WS(6)$	Reference
<i>(not disclosed)</i>	196	-	[G.W. Walker, AMM'50]
Theoretical bound <i>(not proved)</i>	188	554	[J.H. Braun, AMM'50]
SAT	196	572	[Eliahou et al., Computers & Math Applications'12]
Multi-level Tabu-Search	196	574	[Fonlupt et al., EA'11]
SAT <i>(no certificate)</i>	196	575	[Eliahou et al., Computers & Math Applications'12] <i>(revised)</i>

**New:  $WS(6) \geq 581$**

---

## *Successful Key Streamliners:*

{Ordered sets, constrained minimum of each set, partial assignments, sequences of consecutive integers, sequence interleaving}

# Application to the *Weak Schur* problem

## *Successful Key Streamliners:*

{Ordered sets, constrained minimum of each set, partial assignments, sequences of consecutive integers, sequence interleaving}

1 2 4 8 11 22 25 50 63 68 139 149 154 177 182 192 198 393 398 408 413 436 450 455 521 526 540 563 568 578
3 5-7 19 21 23 51-53 64-66 136-138 150-152 179-181 193-195 395-397 409-411 438-440 451-453 523-525 536-538 565-567 579-581
9 10 12-18 20 54-62 140-148 183-191 399-407 441-449 527- 535 569-577
24 26-49 153 155-176 178 412 414-435 437 539 541-562 564
67 69-135 454 456-520 522
196 197 199-392 394

**Partition of [1,581] into 6 weakly sum-free sets, proving  $WS(6) \geq 581$ .**

(le Bras, Gomes,  
Selman AAI 2012)

■ Incrementally obtained, although **not** an example of a **fully constructive procedure** yet, any progress on Schur numbers is quite **significant** given their long history.

# Summary

---

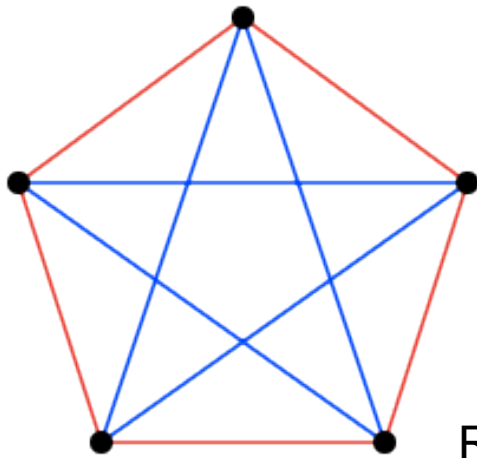


- **Structure discovery** through a general framework that integrates **streamlined** combinatorial search with **human insight and intuition** in an **iterative** approach to discover **efficient constructive procedures**.
- Provides the **first constructive procedures** for the *Spatially-Balanced Latin Square* problem.
- Also, improves the best known **lower bound** for the *Weak Schur Number* problem, and, most recently, for “Graceful graphs.”

■ Extensions **crowd-source** the search for regularities in the solution sets (on Mechanical Turk).

And much, much more ambitiously...

- Ramsey numbers (Erdos), and
- possibly, new types of algorithms.



here

5 nodes

red/blue edges

no blue or red triangle

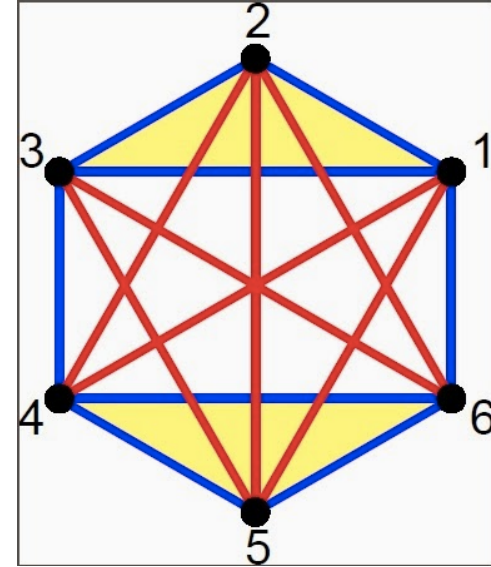
6+ nodes

**always** a blue or red triangle

$R(3,3) = 6$

Ramsey theory about “forced” structure

“unavoidable structure in our universe”



**Erdoes:** *“Imagine an alien force, vastly more powerful than us landing on Earth and demanding the value of  $R(5, 5)$  or they will destroy our planet.” [hmm?]*

$R(3,3) = 6$

$R(4,4) = 18$

$R(5,5) = ?$

*In that case, we should marshal all our computers and all our mathematicians and attempt to find the value. ...*

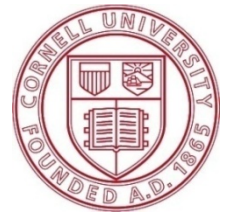
*But suppose, instead, that they asked for  $R(6, 6)$ , then we should attempt to destroy the aliens”.*

Situation is probably not quite as dire. ☺

The key will be to find the hidden structure in the solution space.



# Part II: Crowdsourcing for Backdoor Variables: Materials Discovery Domain



**Briefly!**

---

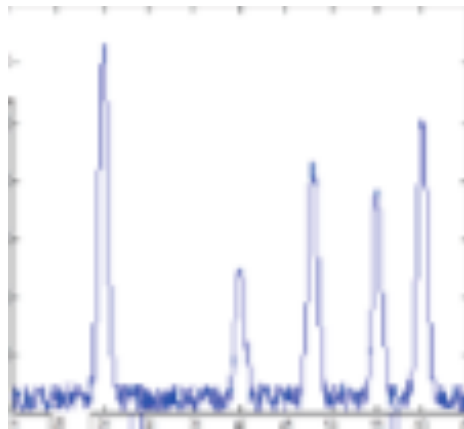


**Warning: Real-world, noisy  
data ahead**

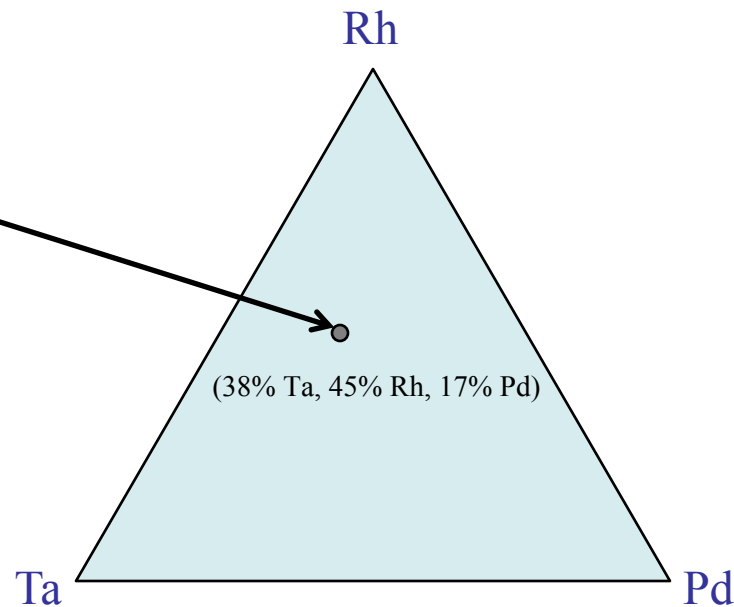


# Motivation

Study of new materials:  
Identifying crystalline structure using **X-Ray Diffraction patterns**



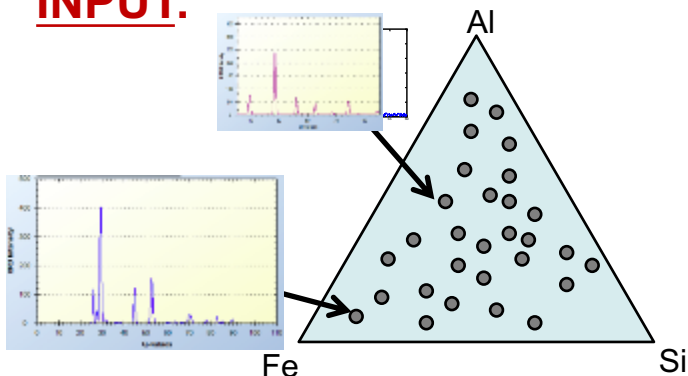
[Source: *Pyrotope*, Sebastien Merkel]



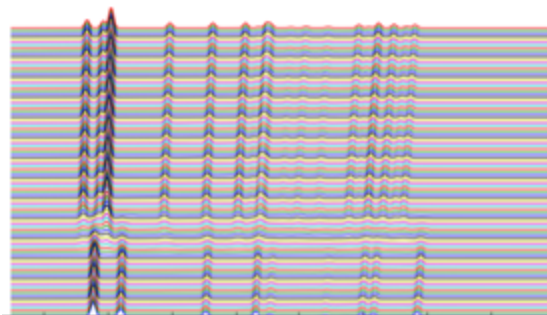
**With Caltech, working  
on screening 1 million  
samples per day.**

# Phase Map Identification Problem

## INPUT:



Collection of X-Ray diffraction Patterns. **Noisy and hidden structure.**

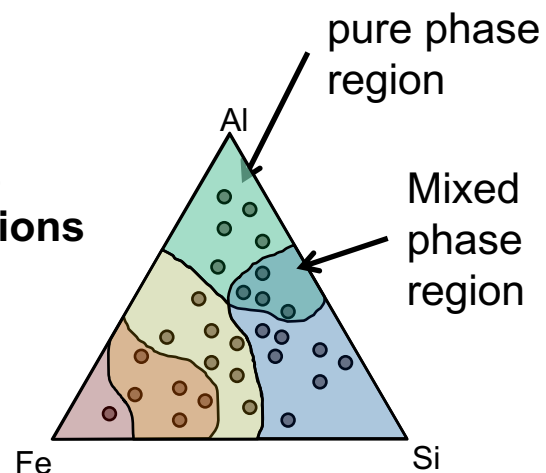


## Additional Physical Requirements:

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

## OUTPUT:

- $m$  phase regions
- $k$  pure regions
  - $m-k$  mixed regions



# Complex Full Constraint Optimization Model

Variables	Description
$p_{ki}$	Normalizing element for phase $k$ in pattern $P_i$
$a_{ki}$	Whether phase $k$ is present in pattern $P_i$
$q_k$	Set of normalized peak locations of phase $B_k$

$$(a_{ki} = 0) \iff (p_{ki} = 0) \quad \forall 1 \leq k \leq K, 1 \leq i \leq n \quad (1)$$

$$1 \leq \sum_{s=1}^K a_{si} \leq M \quad \forall 1 \leq i \leq n \quad (2)$$

$$(p_{ki} = j) \Rightarrow (q_k \subseteq r_{ij}) \quad \forall 1 \leq k \leq K, 1 \leq i \leq n, 1 \leq j \leq |P_i| \quad (3)$$

$$(p_{ki} = j \wedge \sum_{s=1}^K a_{si} = 1) \Rightarrow (r_{ij} \subseteq q_k) \quad \forall 1 \leq k \leq K, 1 \leq i \leq n, 1 \leq j \leq |P_i| \quad (4)$$

$$(p_{ki} = j \wedge p_{k'i} = j' \wedge \sum_{s=1}^K a_{si} = 2) \Rightarrow (\text{member}(r_{ij}[j''], q_k) \vee \text{member}(r_{i'j'}[j''], q_{k'})) \quad \forall 1 \leq k, k' \leq K, 1 \leq i \leq n, 1 \leq j, j', j'' \leq |P_i| \quad (5)$$

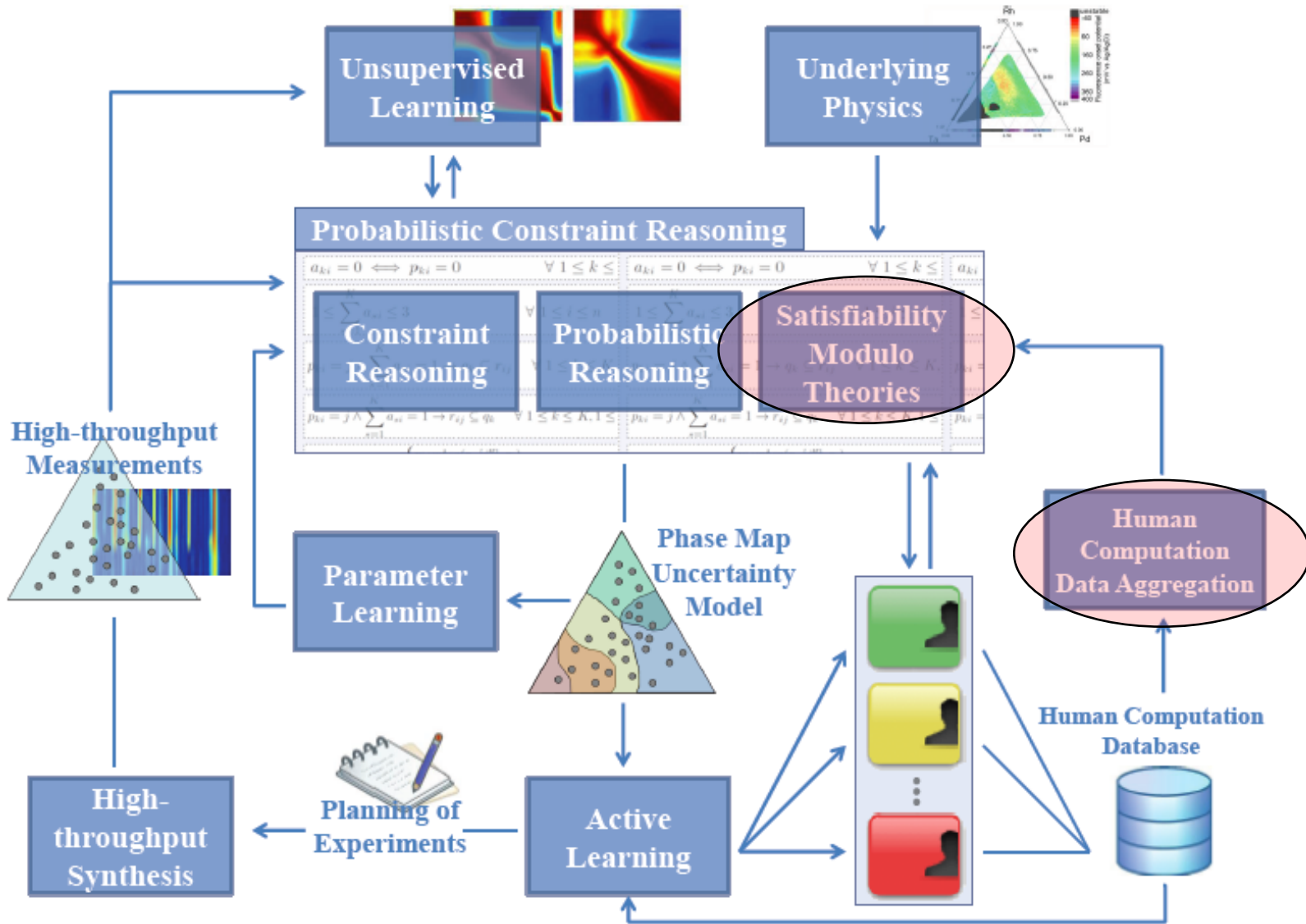
$$(p_{ki} = j) \Rightarrow (p_{k'i} \neq j') \quad \forall 1 \leq k \leq K, (i, j, i', j') \in \Phi \quad (6)$$

$$\Phi = \{(i, j, i', j') \mid \frac{P_i[j]}{P_i[j']} < 1/\delta \vee \frac{P_i[j]}{P_i[j']} > \delta, i < i'\}$$

$$\text{basisPatternConnectivity}(\{a_{ki} \mid 1 \leq i \leq n\}) \quad \forall 1 \leq k \leq K \quad (7)$$

## **Drawback:**

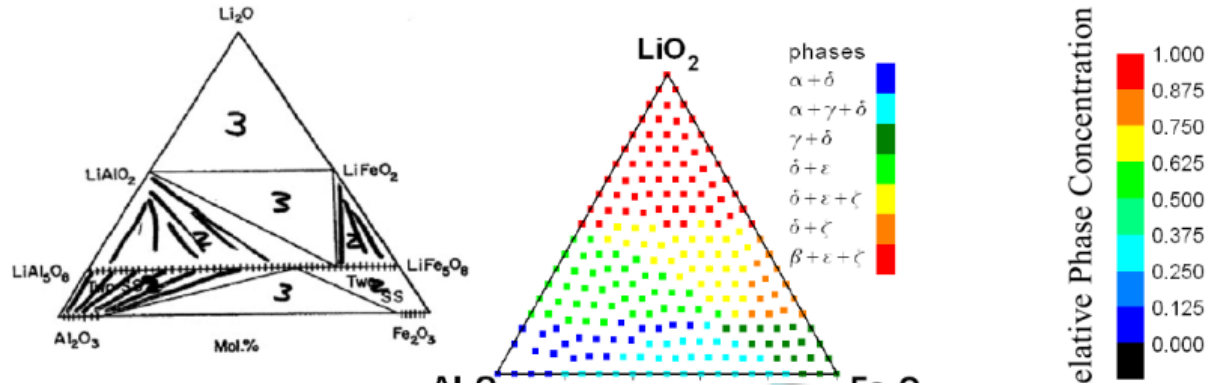
**Does not scale to full-scale realistic instances; poor propagation of experimental noise.**



# SAT Modulo Theory (SMT) Approach: Experimental Validation

Al-Li-Fe instance  
with 6 phases:

Ground truth  
(known)



Previous work



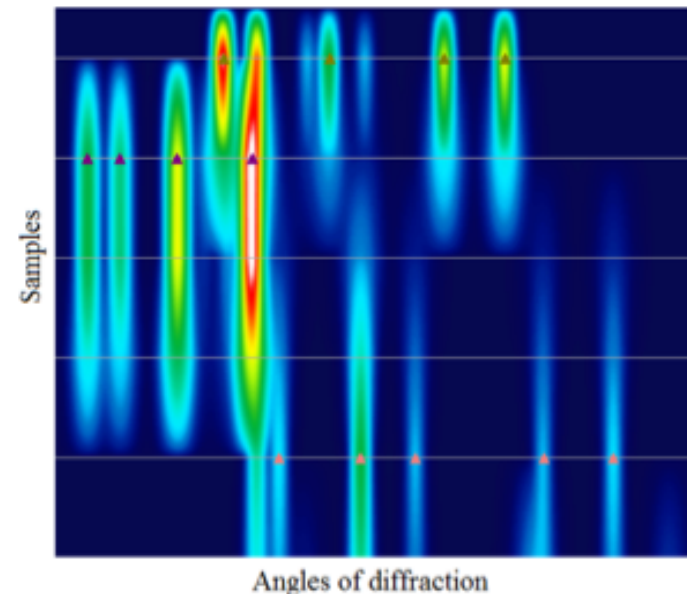
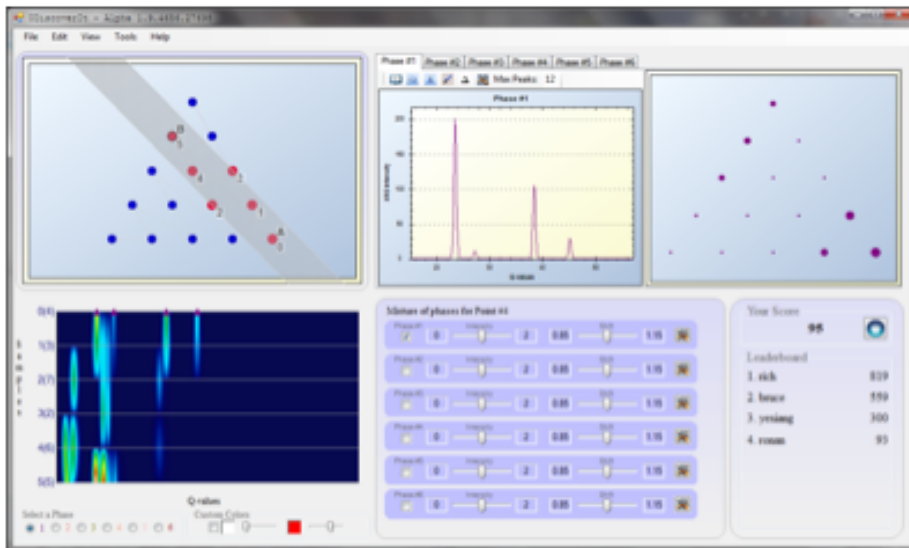
(ML Based approaches)  
violates many  
physical requirements



Our SMT  
approach is much  
more robust



- Can **human input** boost the performance of combinatorial reasoning and optimization methods (SMT)?
  - Human input about potential grouping of diffraction peaks collected on Amazon's Mechanical Turk. (About 500 patterns)
- Provides info on potential so-called **backdoor variables**.





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

Human input can dramatically speed up the performance of combinatorial optimization methods (identification of backdoor variables).

Leverages the complementary strength of **human input**, providing global insights into problem structure, and the **power of combinatorial solvers** to exploit complex constraints.

Note: Human input identifies subsets of peaks that (quite likely) go together.

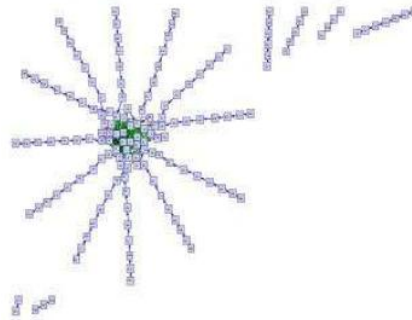
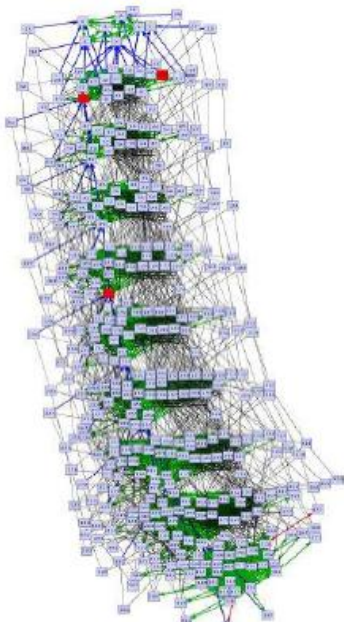


# Observation

**Visual representation** allows human input to be effective in discovering structure and setting backdoor variables in this domain.

**Contrast with earlier work:**

**SAT encodings of planning problems. It was very difficult to give any semantic interpretation to the backdoor variables of the problem instances. Synthetic planning formulas with  $\log(n)$  size backdoors (provably).**



Setting 2 out of 392 variables.

(Hoffmann, Gomes, Selman AIPS-06)

We have discussed how human insights can be combined with reasoning and optimization methods for scientific discovery.

**Finite Mathematics** --- Obtained first procedure for *Spatially Balanced Latin Square* construction.

**Materials Discovery** --- Boosted interpretation of X-ray diffraction patterns with crowdsourced input.