

## ERATO 湊離散構造処理系プロジェクトの 主な研究成果と近況について

### Principal Results and Recent Activities of ERATO MINATO Discrete Structure Manipulation System Project

湊 真一<sup>1\*</sup>

Shin-ichi Minato

**概要** 計算機が扱う問題の多くは、離散構造の処理を基盤としている。近年、論理や集合のような基本データ構造を効率よく処理する「BDD」「ZDD」と呼ばれるデータ構造とアルゴリズムが様々な分野で活用されている。このような技法をベースとして、種々の離散構造を統合的に演算処理する技法を体系化し、分野横断的かつ大規模な実問題を高速に処理する技術基盤を構築することを目標として、「ERATO 湊離散構造処理系プロジェクト」が2009年10月に採択された。今年度までの5年間に、多くの興味深い研究成果が得られている。本講演では、これまでのプロジェクトの主な研究成果と最近の話題、および今後の展望について述べる。

**キーワード** BDD, ZDD, 二分決定グラフ, 離散構造, アルゴリズム, データ構造, 列挙索引化

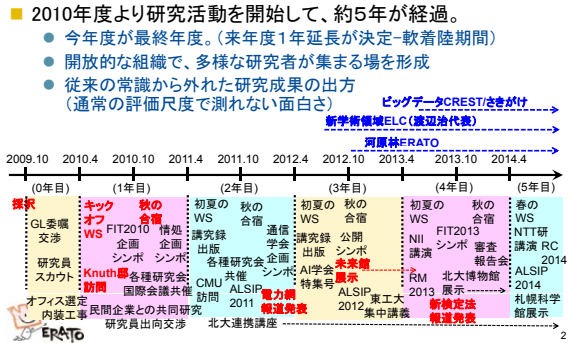
---

1 北海道大学大学院情報科学研究科 / JST ERATO 湊離散構造処理系プロジェクト,  
〒060-0814 札幌市北区北14条西9丁目  
Graduate School of Information Science and Technology, Hokkaido University / JST ERATO MINATO Discrete  
Structure Manipulation System Project,  
North 14 West 9, Sapporo 060-0814, Japan  
\* E-mail address: minato@ist.hokudai.ac.jp

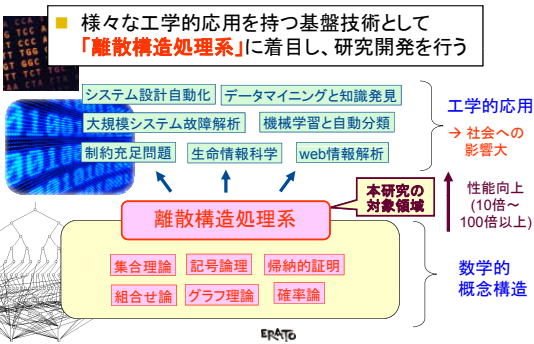
## ERATO 湊離散構造処理系プロジェクトの 主な研究成果と近況について

湊 真一  
北海道大学 / JST ERATO

## 湊ERATOプロジェクトの活動経過

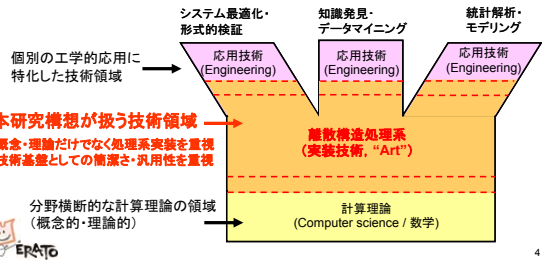


## 本プロジェクトの基本構想

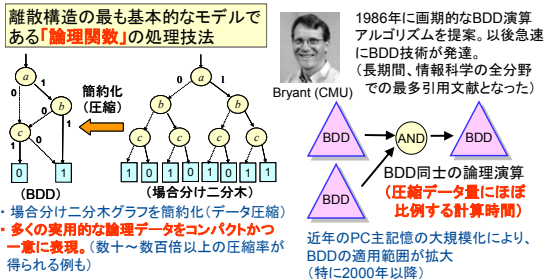


## 本研究構想の対象領域

BDD/ZDD技術の新しい切り口として、様々な離散構造を**統合的に演算処理**する技法を体系化し、**分野横断的かつ大規模な実問題**を高速に処理するための技術基盤を構築する。



## BDD(二分決定グラフ)



## ZDD(ゼロサプレス型BDD)による集合族表現

### 「集合族」を効率的に表現するためのBDDの改良 [Minato1993]

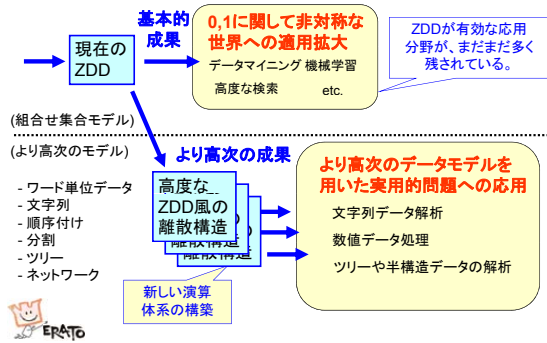
- 集合族: アイテムの組合せを要素とする集合。論理関数と1対1に対応。様々な工学的応用がある。
  - **通常と異なる簡約化規則**を新規に考案。
  - **疎な集合の族を扱う場合に著しい効果**が得られる。(例: 商店の陳列アイテム数に比べて1顧客の購入点数は極めて少ない。)
- 
- 通常のBDD
- ZDDの簡約化

### ZDDはBDDの改良技術として現在、世界的に広く使われている。

- 最近では、データマイニング分野に適用されて、画期的な有効性が示されている。(数百倍のデータ圧縮率・数十倍の処理高速化)
- 他にも応用例は増えつつある。

Dataset name	'manufroom'	CPU time (sec.)	
Min.support	#Items	by ZDD	by BDD
1,000	123,287	0.50	0.64
500	1,442,504	1.32	3.29
300	5,259,786	2.25	9.96
200	18,094,822	3.21	31.63
100	66,076,586	5.06	114.21
70	153,336,056	7.16	277.15
50	198,169,886	8.17	357.27

## 本研究プロジェクトの技術面のポイント



7

## 研究実施場所とメンバ



8

## 初夏のERATOワークショップ / 秋のERATO合宿

- 初夏のERATOワークショップ (毎年5~6月に開催)
  - 北大ERATOオフィスで開催(1泊2日)
  - 正規メンバを含む40~60名の研究者が参加。
  - 基調講演3~4件とポスター発表30~40件
- 秋のERATO合宿 (毎年10~11月に開催)
  - 北海道内の旅館/ホテルを半貸切で会合(2泊3日)(北広島、支笏湖、夕張、登別)
  - 正規メンバを含む40~50名の研究者が参加
  - 講演(Q&Aつき)約25件とポスター発表 約15件。
  - 深夜~未明まで続く熱いディスカッション

- 発表者の質・量ともに素晴らしい  
- 未発表の新鮮なアイデアを真剣に議論  
- 新しい共同研究に発展することも



9

## ERATO 세미나と講究録

- ERATO 세미나: ほぼ毎週・年間35~40回実施。これまでに**150回を突破**。
  - 国内外の第一線の研究者による講演と意見交換
  - セミナーの後に有志で懇親会を実施
- セミナー講演およびワークショップ・合宿での講演資料を再編集して「講究録」として毎年出版している。CD-ROMは関係者に無償配布。
  - 講演者の研究分野が幅広く、新しい研究の種が含まれている。



10

## 本プロジェクトの主な研究成果

## 特筆すべき成果

- フロンティア法と電力網への応用
  - 社会経済への波及効果大
  - 電力網の網羅的解析を世界で初めて実現
- ビッグデータから新たな科学的発見をもたらす統計手法
  - 実験系のあらゆる学問への波及効果大
  - PNASIに掲載。将来ものすごく引用される可能性あり
- 「フカシギの数え方」の展示と数え上げ世界記録
  - 青少年・一般市民への波及効果大
  - 展示来場者23万人、YouTube155万ビュー



11



12

## 電力網への応用(2012年2月プレスリリース)

- 林 泰弘 教授(早稲田大)との共同研究
  - 電力系スマートグリッド業界のリーダ的存在(経産省スマートハウス標準化検討会座長、他多数の要職)
  - 電力網最適化の研究で1990年代より濃と協力関係

- 大震災後、より緊急性の高い研究課題に
  - 今後長期的に不足する電力を自然エネルギーで補うために必須の電力網解析・制御技術を支援

情報科学の研究者集団として我が国の苦境を克服するためできる限り貢献したい。

→ ERATOプロジェクトでの取り組みを加速



## フロンティア法と電力網への応用

「超高速アルゴリズム」にできること

- 配電網の各スイッチのON/OFFを、ZDDの値を元に決定づけ
- 制約を違えずかどうかを終端節点の記号に対応づけ
- 制約条件を満たす (OK)
- 制約条件を満たさない (NG)

ZDD技法で任組して超高速に処理

大震災後に重要性を増した電力網の解析技法の高速化・大規模化・高信頼化

## 標準的な電力網モデルへの適用結果

- ・ 標準的な電力網モデル(スイッチ468個)でグラフ的制約と電氣的制約を共に満たす解を全て求め、圧縮して表現することに成功
- ・ 電力の損失を最小にする組合せを求めることに成功

圧縮データ: 約110万個 (779MB) 実行時間: 数十分  
 解の個数: 約10<sup>63</sup> (2136那由他8201阿僧祇3834恒河沙8532億9116載8261正2214彌6049満560穰9817穉8392垓4438京5235兆3981億8952万1540) 通り

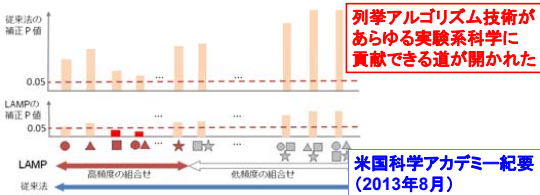
- ・ 普通の方法で探したのでは絶対終わらない数だが、ZDDで圧縮してまとめて、正確に全て探し出している
- ・ 計算結果は計算機メモリに索引化されて保存

Takeru Inoue, Keiji Takano, Takayuki Watanabe, Jun Kawahara, Ryo Yoshinaka, Akihiro Kishimoto, Koji Tsuda, Shin-ichi Minato, and Yasuhiro Hayashi: "Distribution Loss Minimization With Guaranteed Error Bound," *IEEE Trans. Smart Grid*, Vol. 5, No. 1, pp. 102-111, Jan. 2014.

## 特筆すべき成果

- フロンティア法と電力網への応用
  - 社会経済への波及効果大  
電力網の網羅的解析を世界で初めて実現
- ビッグデータから新たな科学的発見をもたらす統計手法
  - 実験系のあらゆる学問への波及効果大  
PNASIに掲載。将来ものすごく引用される可能性あり
- 「フカシギの数え方」の展示と数え上げ世界記録
  - 青少年・一般市民への波及効果大  
展示来場者23万人、YouTube155万ビュー

## ビッグデータから新たな科学的発見をもたらす統計手法



### Statistical significance of combinatorial regulations

Aika Terada<sup>1,2,3</sup>, Mariko Okada-Hatakeyama<sup>4</sup>, Koji Tsuda<sup>4,5,6</sup>, and Jun Sese<sup>1,2,3</sup>

<sup>1</sup>Department of Computer Science and <sup>2</sup>Education Academy of Computational Life Sciences, Tokyo Institute of Technology, Meguro-ku, Tokyo 152-8550, Japan; <sup>3</sup>Minato Discrete Structure Manipulation System Project, Exploratory Research for Advanced Technology, Japan Science and Technology Agency, Sapporo, Hokkaido 050-8584, Japan; <sup>4</sup>Laboratory for Integrated Cellular Systems, RIKEN Center for Integrated Medical Sciences (RIMS-ICMS), Yokohama, Kanagawa 220-0045, Japan; and <sup>5</sup>Computational Biology Research Center, National Institute of Advanced Industrial Science and Technology, Koto-ku, Tokyo 135-0064, Japan

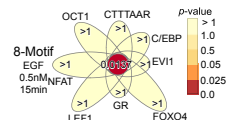
Edited by Wing Hung Wong, Stanford University, Stanford, CA, and approved July 3, 2013 (received for review February 4, 2013)

More than three transcription factors often work together to enable cells to respond to various signals. The detection of combinatorial regulation by multiple transcription factors, however, is not only deliberately excluding such tests. Here, we propose an efficient branch-and-bound algorithm, called the "limitless arity multiple testing procedure" (LAMP). LAMP counts the exact number of

## LAMP (Limitless Arity Multiple testing Procedure)

- ERATO東京グループの成果(寺田・瀬々・津田)
- 科学雑誌では、検定値(P値)が十分小さくないと、発見として認められない
- 多重検定補正(ボンフェローニ法): n個の要因候補から発見する場合、0.05/nをP値の閾値として用いる
- 組合せ要因の数は膨大 ⇒ 発見は非常に困難
- 列挙アルゴリズムで出現頻度の低い組合せ要因を除く: 発見力の大幅向上

これまで見過ごされてきた、乳がん細胞の分化に関与する転写因子の組合せを発見!



## LAMP法のキーアイデア

- 「出現頻度が非常に低い組合せパターンは、決して有意にならないので、仮説集合から除外しても判定に影響しない。」
  - あるパタンの出現頻度を $\lambda$ とすると、それが全部正例に入るという最も極端な偏りを起こす確率が  $\frac{p_{C_A}}{m \cdot C_A}$  なので、これが偶然のデータ偏りの確率値(P値)の下限となる。(以降、これを $f(\lambda)$ と書く)
- P値の閾値が $f(\lambda)$ より小さければ、出現頻度 $\lambda$ 以下のパターンは決して有意にならないので、仮説集合から除外できる。
  - 出現頻度 $\lambda$ 以上のパタンの総数 $k_\lambda$ を数え上げてこれを仮説集合とすると、補正した閾値は、 $\alpha / k_\lambda$ となる。
  - $\alpha / k_\lambda < f(\lambda)$  となるような、最大の頻度 $\lambda_{\max}$ を見つけ出せばよい。

19

## LCMアルゴリズムに基づくLAMPの計算

- LCM: NII宇野先生が開発した頻出パターンマイニング手法
  - 世界最高速の実装
- LAMP法のホームページでプログラムが公開されている。
  - 内部で宇野先生のLCM実装を呼び出している。
- 結構、計算時間がかかっていて、少し大きなデータになると計算に何日もかかる、という問題を抱えていた。
  - 最近、新しいアルゴリズムの開発に成功 (ECML/PKDD2014に採択)
  - 従来、頻度1きざみで何回もLCMを呼び出して探索していた。
    - 1回の深さ優先探索でLAMPの閾値を求めることが可能に

20

## 特筆すべき成果

- フロンティア法と電力網への応用
  - 社会経済への波及効果大
  - 電力網の網羅的解析を世界で初めて実現
- ビッグデータから新たな科学的発見をもたらす統計手法
  - 実験系のあらゆる学問への波及効果大
  - PNASIに掲載。将来ものすごく引用される可能性あり
- 「フカシギの数え方」の展示と数え上げ世界記録
  - 青少年・一般市民への波及効果大
  - 展示来場者23万人、YouTube155万ビュー



21

## 「フカシギの数え方」の展示と数え上げ世界記録

YouTube 未来館チャンネルで 1,557,869 ビュー (2014年8月17日現在)

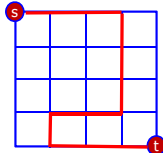
組み合わせの数え方

情報科学の重要性を 小中高生や一般市民に 強く印象付けることに成功

22

## “self-avoiding walks”の数え上げ

- 最短経路の数え上げは簡単 (→  $2n C_n$ ; 高校で習う問題)
- 最短でない経路を許すと突然難しくなる。(計算式や漸化式は見つかっていない)
  - おねえさんが25年かかった結果がアニメの中に表示されているため、「計算式教えて」というコメント多数。
  - 残念ながら計算式は知られていない。効率良く数え上げるしかない。



23

## 「おねえさんの問題」の世界記録

- ERATO研究員の岩下氏が世界記録を更新(2013年12月)
  - n=26までの数え上げに成功
  - ノルウェーの大学が記録したn=24(2013年2月)を2段階上回る。
  - 様々な工夫を凝らした末に、主記憶2TBの計算機2台を約1週間占有して計算。

ERATO logo

24

情報処理学会誌 解説記事(2013年11月号)

◆解説◆  
「おねえさんの問題」の最先端  
—YouTube 動画と世界記録—  
渡 真一 (北海道大学大学院情報科学研究科 / JST ERATO 演算数値最適化基盤プロジェクト)

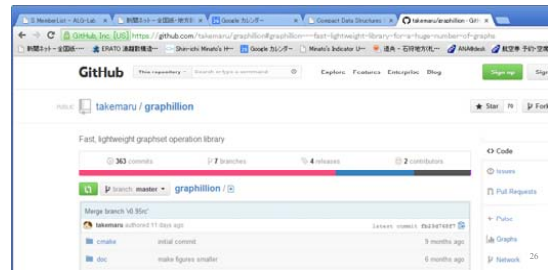
「おねえさんの問題」とは  
2012年8月1日より2013年4月15日までの約  
6ヵ月間、東京大学理学院の大学院生による  
読者アンケート評価  
2013年の全記事中で第1位

図-1 「オカンキおねえさん」の動画のスクリーンショット  
の一例として、組合せ爆発のすごさとアルゴリズム  
あるが、最先端のアルゴリズム技術を使えば、同じ

25

Open software: "Graphillion.org"

- Toolbox for ZDD-based graph enumeration.  
– Easy interface using Python graph library.



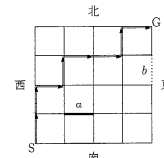
Graphillion チュートリアル動画



27

今年の大学入試問題 (2014年度 北大理系 数学)

- 4 図のような格子状の道路がある。S 地点を出発して、東または北に進んで G 地点に到達する経路を考える。ただし太い実線で描かれた区間  $a$  を通り抜けるのに 1 分、点線で描かれた区間  $b$  を通り抜けるのに 8 分、それ以外の各区間を通り抜けるのに 2 分かるとする。たとえば、図の矢印に沿った経路では S を出発し G に到達するまでに 16 分かかる。



- (1)  $a$  を通り抜ける経路は何通りあるか。
- (2)  $a$  を通り抜けずに  $b$  を通り抜ける経路は何通りあるか。
- (3) すべての経路から任意に 1 つ選んだとき、S 地点から G 地点に到達するのにかかる時間の期待値を求めよ。

ERATO

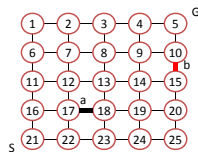
28

```
>>> from graphillion import GraphSet
>>> import graphillion.tutorial as tl
>>> GraphSet.set_universe(tl.grid(4, 4))
>>> S, G, a, b, L = 21, 5, (17,18), (18,15), 8
>>> P = GraphSet.paths(S, G).len(L)

>>> '(1)', len(P.including(a))
(1) 28

>>> '(2)', len(P.excluding(a).including(b))
(2) 9

>>> def E(P): # 集合 P のパスを通る時間の期待値
>>> L = 0.
>>> for p in P:
>>>     for e in p:
>>>         if e == a: L += 1
>>>         elif e == b: L += 8
>>>         else: L += 2
>>>     return L / len(P)
>>> '(3)', E(P)
(3) 17.0
```



Graphillion で解いた例  
by 井上 武さん  
(元ERATO 研究員-現在 NTT)

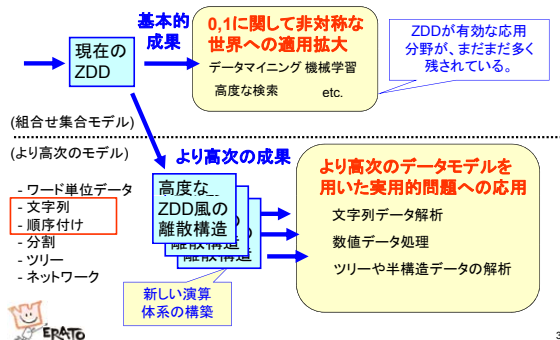
29

遠回りを許して同じところを2度通らない場合でも解ける。  
(こちらの方が Graphillion の記述は簡単)

```
>>> P = GraphSet.paths(S, G)
>>> '(1)', len(P.including(a))
(1) 3322
>>> '(2)', len(P.excluding(a).including(b))
(2) 2698
>>> '(3)', E(P)
(3) 37.4081296992
```

30

## 本研究プロジェクトの技術面のポイント



31

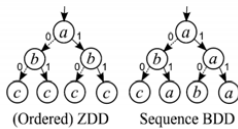
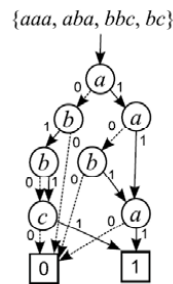
## 系列集合(Sets of Sequences)

- Sets of combinations:
  - Don't consider order and duplication of items
  - "abcc" and "bca" are the same.
- Sets of sequences:
  - Distinguishes all finite sequences.
  - $\varnothing, \{\lambda\}, \{ab, aba, bbc\}, \{a, aa, aaa, aaaa\}, \text{etc.}$
  - Here we exclude infinite sets such as  $a^*$ .
- So many real-life applications.
  - Text search and indexing
  - Web (html/xml) data mining
  - Bio informatics

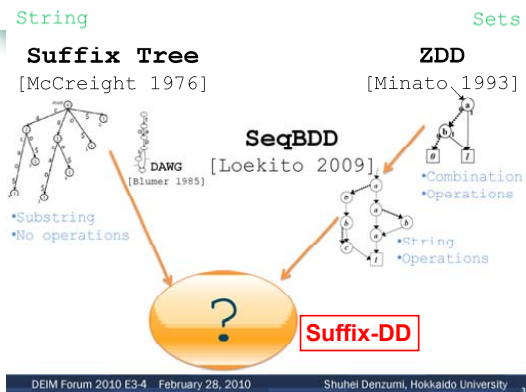
32

## Sequence BDD (SeqBDD)

- Loekito, Bailey, and Pei (2009)
  - Same as ZDD reduction rule.
  - Only 0-edges** keep variable ordering.
  - 1-edges has no restriction.
  - Still unique representation for a given set of sequences.
  - Each path from root to 1-terminal corresponds to a sequence.



33

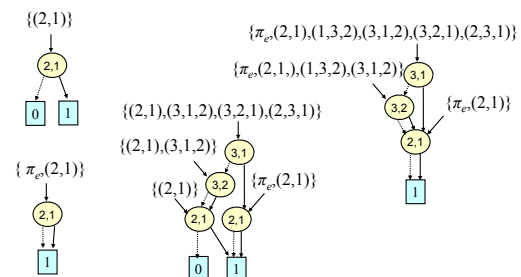


## 順列集合の処理

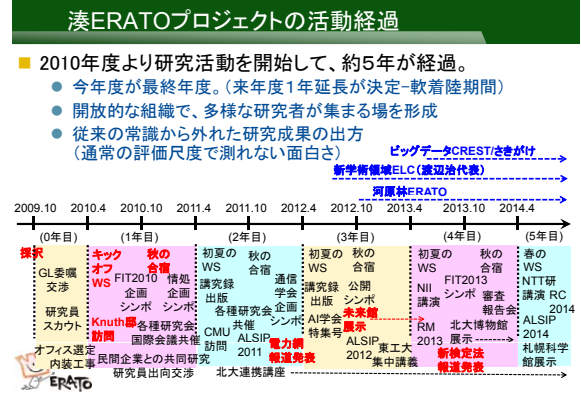
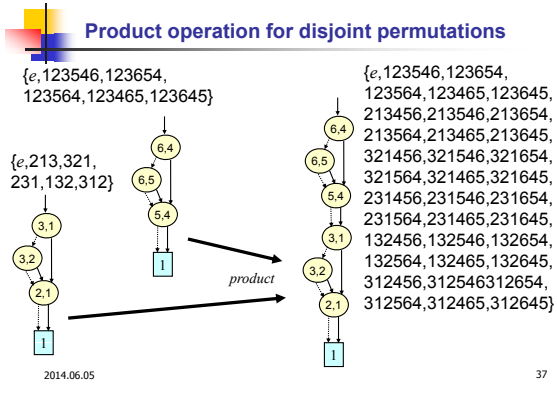
- ルービックキューブ:
  - $P = \{\pi\}$  基本操作による順列
  - $P$  は18通り (= 3通り × 6面) の順列の集合
  - 直積演算  $P \times P$  は、2回の基本操作で得られるパタンの集合を表す。
  - $P^{20}$  で原理的には全ての可能なパタンの集合が得られる。(実際は場合の数が多すぎてメモリあふれが起きる)
- 15パズル・トランプ等も順列集合で表現される
  - 最適な操作系列の解析
- あみだくじ・ソーティングネットワーク
  - 1対1のマッチング問題に適用可能
  - 順列は単射写像に対応。群論と深い関係がある。
- 情報無損失符号の設計・暗号理論への応用
  - Reversible logic設計問題への応用 (量子計算とも関連)

35

## Permutation Decision Diagram ( $\pi$ DD)



36



- ### 今年度～来年度の活動予定
- ERATOワークショップ(北大 4/17～18) (招待講演:L. Simon 先生)
  - ALSIP2014主催(台湾 5/13) (PAKDD併設ワークショップ)
  - RC2014共催(京都 7/10～11) (可逆計算・量子計算)
  - 通信学会IN研究会共催(北大 7/17～18)
  - ERATO合宿(礼文島 9/7～10)
  - プロジェクト最終報告会(東京 1月下旬)
  - WALCOM2015基調講演(バングラデシュ 2/26～28)
  - プロジェクト研究成果のpublication
    - 電子情報通信学会誌(2014年12月号)小特集(約30ページ) 「BDD/ZDDを基盤とする離散構造処理系とその応用」
    - 専門書執筆
  - 特別重点期間(2015年4月～2016年3月)
    - 年間予算規模1/4程度、雇用研究員のソフトランディング

- ### おわりに
- 情報科学分野のERATO組織をいかに構成するか
    - 研究する「人」と集まる「場」がすべて
    - ピラミッド型ではなく、遊牧的なオープンな組織を構成
    - 研究テーマは一貫性を持たせた。(離散構造・BDD/ZDDを軸に)
    - 非実験系の大型プロジェクトの進め方の1つの典型例を示せたのではないか
  - 「Art層」のプロジェクトは、成果の定量的評価は難しい
    - 最初のきっかけはZDDがKnuth本に取り上げられたこと → そもそもインパクト計測不能
    - 「Art層」の技術分野は、論文件数が少なくても大事な研究がある
    - 質的に突き抜けたERATOらしい成果をアピール
  - 計画してできるようなものでない成果:「芽」を逃さず具現化
    - 研究の方向性、オープンな協力体制、機動的リーダーシップの3点