

# Algebraic Methods in ML

Symposium and workshop



Jason Morton  
Stanford, Math



Guy Lebanon  
CSE, Georgia Tech



Risi Kondor  
Gatsby Unit, UCL

# Themes:

- Algebraic statistics
- Non-commutative harmonic analysis
- Combinatorial problems, e.g., ranking

13.30	<b>Risi Kondor:</b> Non-commutative harmonic analysis
14.05	<b>Guy Lebanon:</b> Modeling distributions on permutations and partial ranking
14.40	<b>Jason Morton:</b> Algebraic models for multilinear dependence
15.15	coffee break
15.25	<b>Yanxi Liu:</b> Symmetry Group-based Learning for Regularity Discovery from Real World Patterns
16.00	<b>Marina Meila:</b> Estimation and model selection in stagewise ranking: a representation story

# Workshop session 1 (Callaghan room @ Westin)



7.30	<b>Stephen E. Fienberg:</b> Algebraic statistics for random graph models: Markov bases and their uses
8.05	<b>Adrian Dobra:</b> Algebraic statistics and contingency tables
8.40	<b>Keisuke Yamazaki:</b> Toric Modification on Mixture Models
9.15	coffee break
9.25	<b>Vincent Auvray:</b> Learning Parameters in Discrete Naive Bayes Models by Computing Fibers of the Parametrization map
10.00	<b>Paul von Bunau:</b> Stationary Subspace Analysis


# Workshop session 1 (Callaghan room @ Westin)

15.30	<b>Doru Balcan:</b> Alternatives to the Discrete Fourier Transform
16.05	<b>Lek-Heng Lim:</b> Graph Helmholtzian and rank learning
16.40	<b>Xiaoye Jiang:</b> Identity Management On Homogeneous spaces
17.15	coffee break
17.25	<b>Jonathan Huang:</b> Exploiting Probabilistic Independence for Permutations
17.55	<b>Tiberio Caetano:</b> Consistent structured estimation for weighted bipartite matching

# Non-commutative harmonic analysis

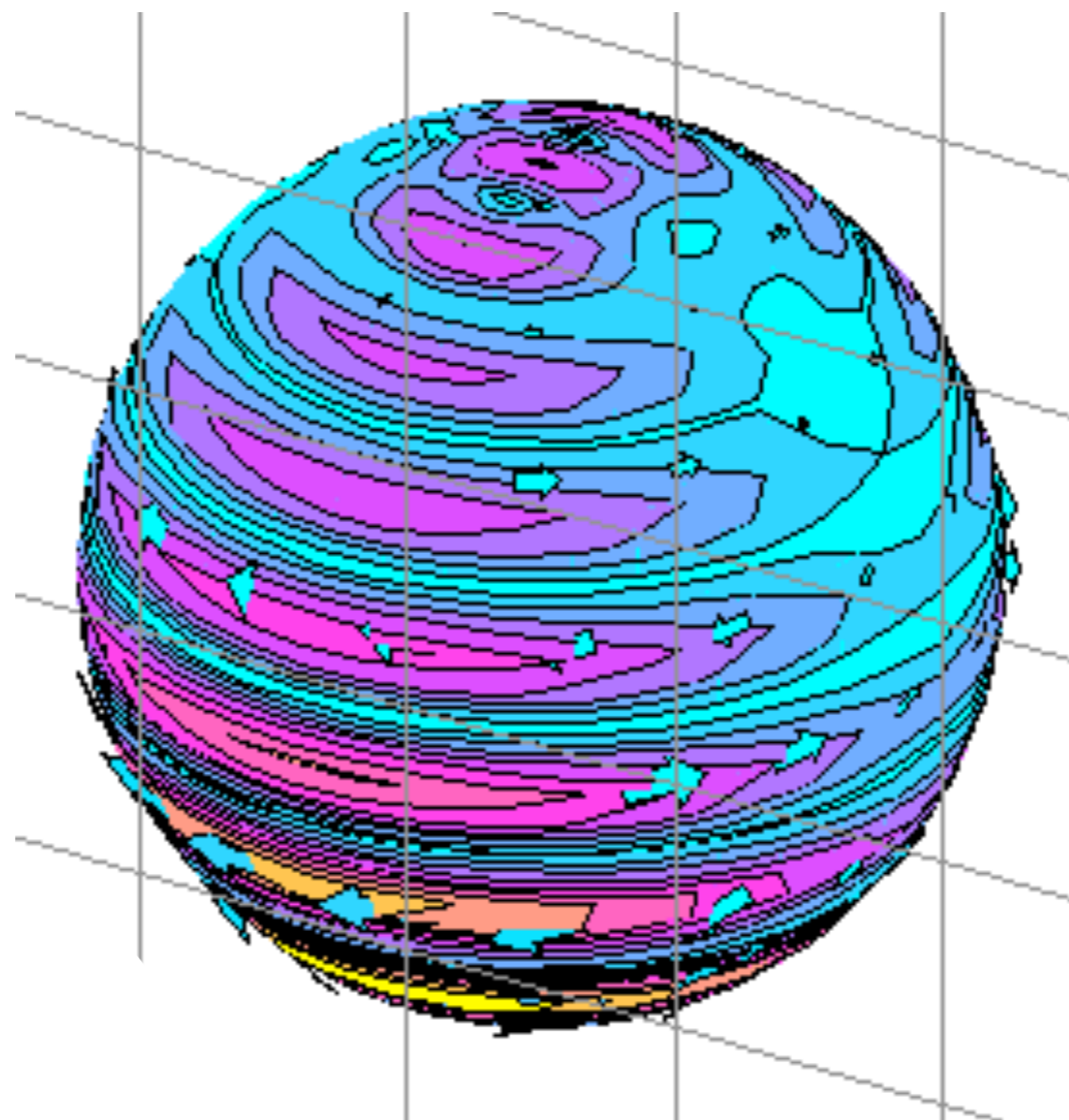
Risi Kondor


$$f(x) = \int e^{2\pi i k x} \hat{f}(k) dk$$


$$\hat{f}(k) = \int e^{-2\pi i k x} f(x) dx$$




$$f : S^2 \rightarrow \mathbb{R}$$



$$[1, 2, 3, 4, 5] \mapsto c_1$$

$$[1, 2, 3, 5, 4] \mapsto c_2$$

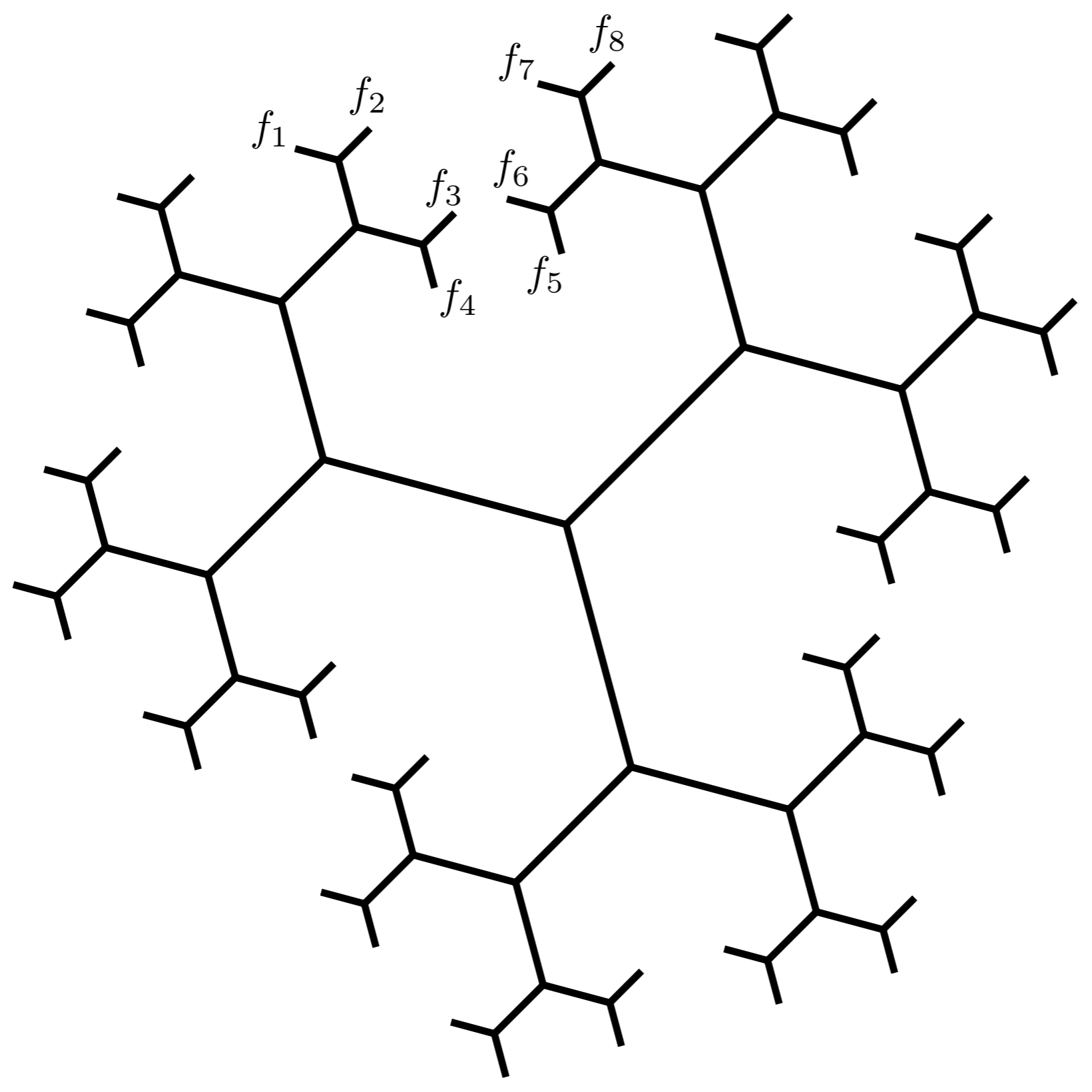
$$[1, 2, 5, 3, 4] \mapsto c_3$$

$$[1, 5, 2, 3, 4] \mapsto c_4$$

$$[5, 1, 2, 3, 4] \mapsto c_5$$

$\vdots \quad \vdots \quad \vdots$

$$[5, 4, 3, 2, 1] \mapsto c_{5!}$$



$$\sqrt{\heartsuit} = ?$$

$$\cos \heartsuit = ?$$

$$\frac{d}{dx} \heartsuit = ?$$

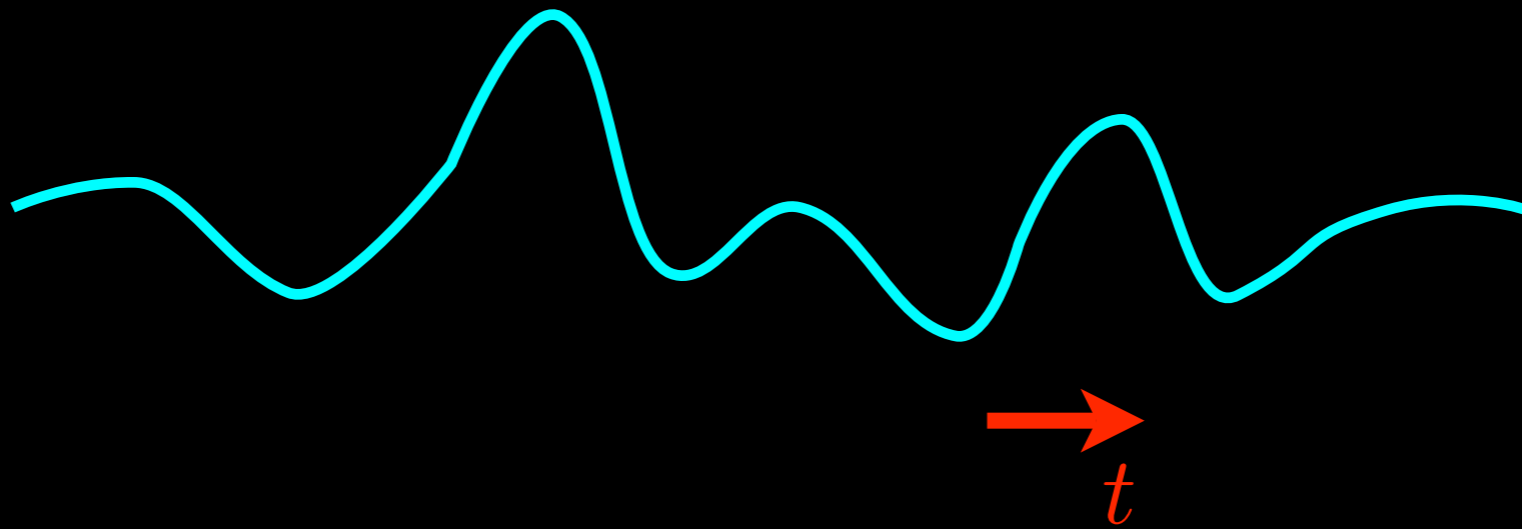
$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \heartsuit = ?$$

$$F\{\heartsuit\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(t) e^{it\heartsuit} dt = ?$$

My normal approach  
is useless here.

$$\hat{f}(k) = \int e^{-i2\pi kx} f(x) dx$$

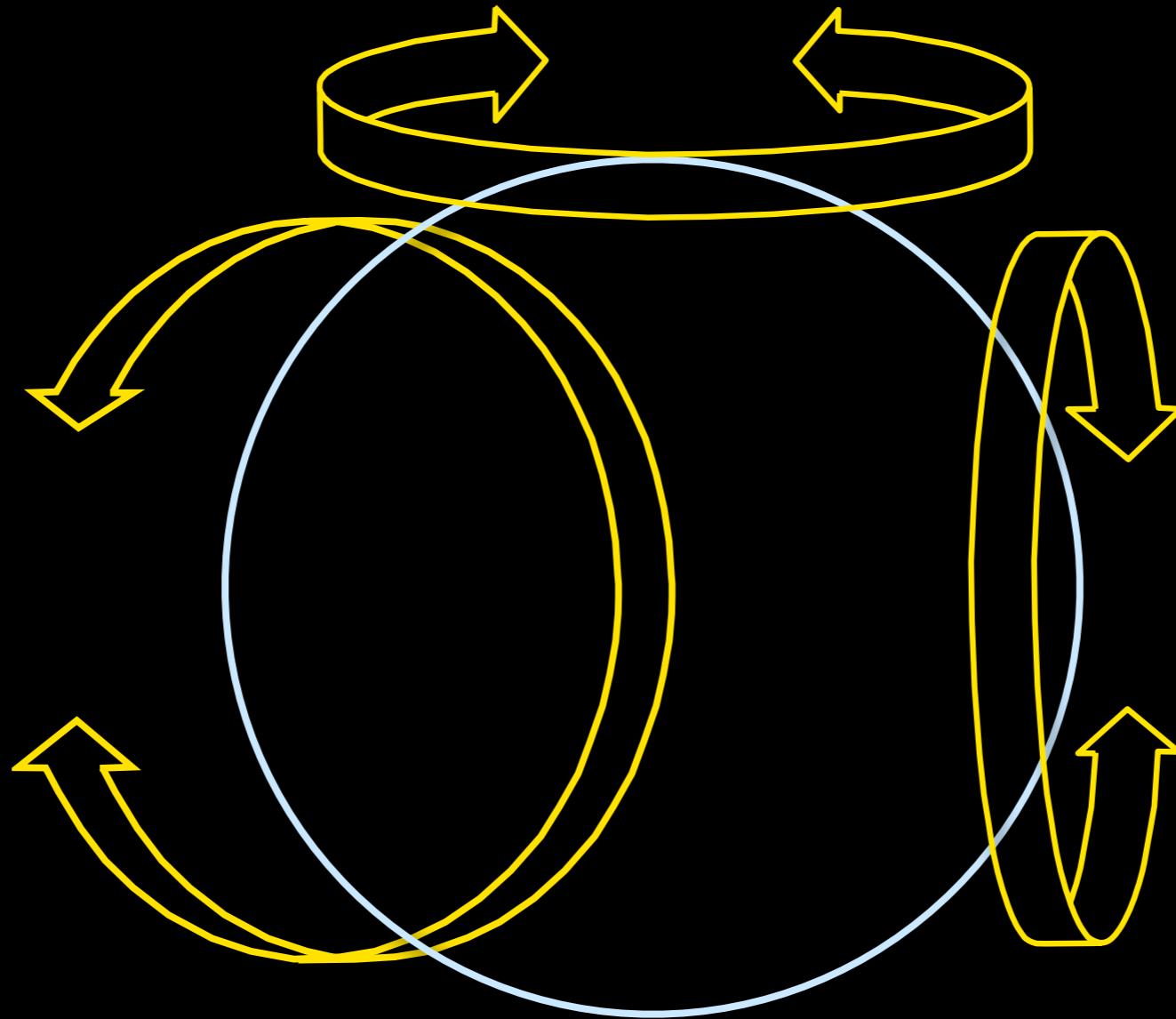
$$f^t(x) = f(x - t)$$



$$\hat{f}^t(k) = e^{-i2\pi kt} \hat{f}(k)$$

$$f^R(x) = f(R^{-1}x)$$

$$R \in \text{SO}(3)$$



In general,  $f: \mathcal{X} \rightarrow \mathbb{R}$  and  $G$  is a **group** acting on  $\mathcal{X}$ .

$$g_1 g_2 \in G \quad \forall g_1, g_2 \in G$$

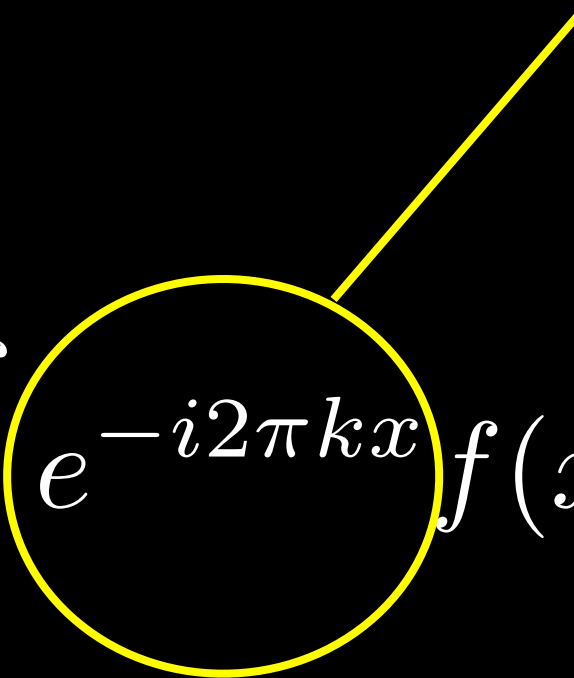
$$g_3 (g_2 g_1) = (g_3 g_2) g_1 \quad \forall g_1, g_2, g_3 \in G$$

$$eg = ge = g \quad \text{for some } e \in G$$

$$g^{-1}g = gg^{-1} = e \quad \text{for some } g^{-1} \in G$$



$$e^{-i2\pi kx_1} e^{-i2\pi kx_2} = e^{-i2\pi k(x_1+x_2)}$$

$$\hat{f}(k) = \int e^{-i2\pi kx} f(x) dx$$


$$\rho(x_2)\rho(x_1) = \rho(x_2x_1)$$

$$\hat{f}(\rho) = \sum_{x \in G} f(x) \rho(x)$$


$\rho: G \rightarrow \mathbb{C}^{d \times d}$  is called a **representation**

# Irreducible Representations of $SO(3)$

$$\rho_0(R) \quad \square \quad [\rho_l(\theta, \phi, \psi)]_{m,m'} = e^{-im'\psi} Y_l^m(\theta, \phi)$$

$$\rho_1(R) \quad \square$$

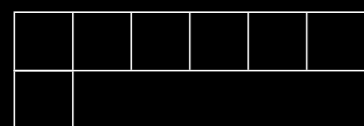
$$\rho_2(R) \quad \square$$

$$\rho_3(R) \quad \square$$

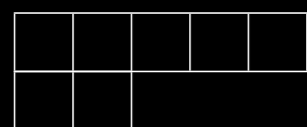
# Irreducible representations of $S_n$



$$d = 1$$



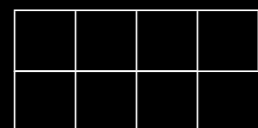
$$d = n - 1$$



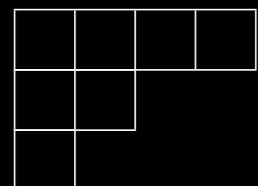
$$d = n(n - 3)/2$$



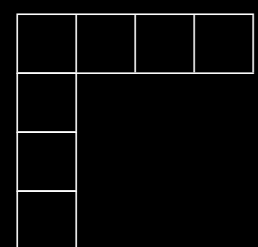
$$d = (n - 1)(n - 2)/2$$



$$d = n(n - 1)(n - 5)/6$$



$$d = n(n - 2)(n - 4)/3$$

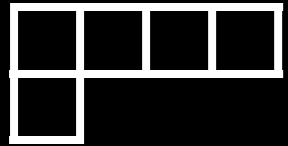


$$d = (n - 1)(n - 2)(n - 3)/6$$

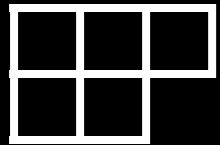
$$\sigma = (2, 1, 3, 4, 5) \in \mathfrak{S}_5$$



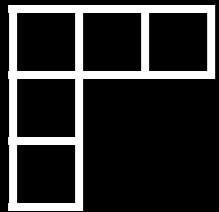
$$\rho_1(\sigma) = (1)$$



$$\rho_2(\sigma) = \begin{pmatrix} -0.5 & -0.289 & -0.204 & -0.791 \\ 0.866 & -0.167 & -0.118 & -0.456 \\ 0 & 0.943 & -0.0833 & -0.323 \\ 0 & 0 & 0.968 & -0.25 \end{pmatrix}$$



$$\rho_3(\sigma) = \begin{pmatrix} 0.25 & -0.433 & 0.433 & -0.75 & 0 \\ -0.433 & -0.25 & -0.75 & -0.433 & 0 \\ -0.433 & -0.25 & 0.25 & 0.144 & -0.816 \\ 0.75 & -0.144 & -0.433 & 0.0833 & -0.471 \\ 0 & 0.816 & 0 & -0.471 & -0.333 \end{pmatrix}$$



$$\rho_4(\sigma) = \begin{pmatrix} 0.333 & 0.236 & 0 & 0.913 & 0 & 0 \\ -0.471 & 0.0417 & 0.217 & 0.161 & 0.839 & 0 \\ 0.816 & -0.0722 & 0.125 & -0.28 & 0.484 & 0 \\ 0 & -0.484 & -0.28 & 0.125 & 0.0722 & 0.816 \\ 0 & 0.839 & -0.161 & -0.217 & 0.0417 & 0.471 \\ 0 & 0 & 0.913 & 0 & -0.236 & 0.333 \end{pmatrix}$$

$$\widehat{f}(\rho) = \sum_{x \in G} f(x) \rho(x) \quad f(x) = \frac{1}{|G|} \sum_{\rho \in \mathcal{R}} d_\rho \operatorname{tr} \left[ \widehat{f}(\rho) \rho(x^{-1}) \right]$$

1. **Linearity:**  $\widehat{f + g} = \widehat{f} + \widehat{g}$

2. **Unitarity:**  $\langle f, g \rangle = \langle \widehat{f}, \widehat{g} \rangle$

3. **Left-translation:**  $\widehat{f^z}(\rho) = \rho(z) \widehat{f}(\rho)$

4. **Convolution:**  $\widehat{f * g}(\rho) = \widehat{f}(\rho) \widehat{g}(\rho)$

5. The individual components correspond to different levels of smoothness.

Group algebra:  $\{e_x\}_{x \in G}$   $e_x \cdot e_y := e_{xy}$

$$\mathcal{F}: \mathbb{C}G \xrightarrow{\text{isom.}} \bigoplus_{\rho \in \mathcal{R}} \text{GL}(d_\rho)$$



“Group Representations in  
Probability and Statistics”  
IMS, 1988

Persi Diaconis



**Multi-object tracking with representations of the symmetric group [K., Howard, Jebara, AISTATS 2007]**

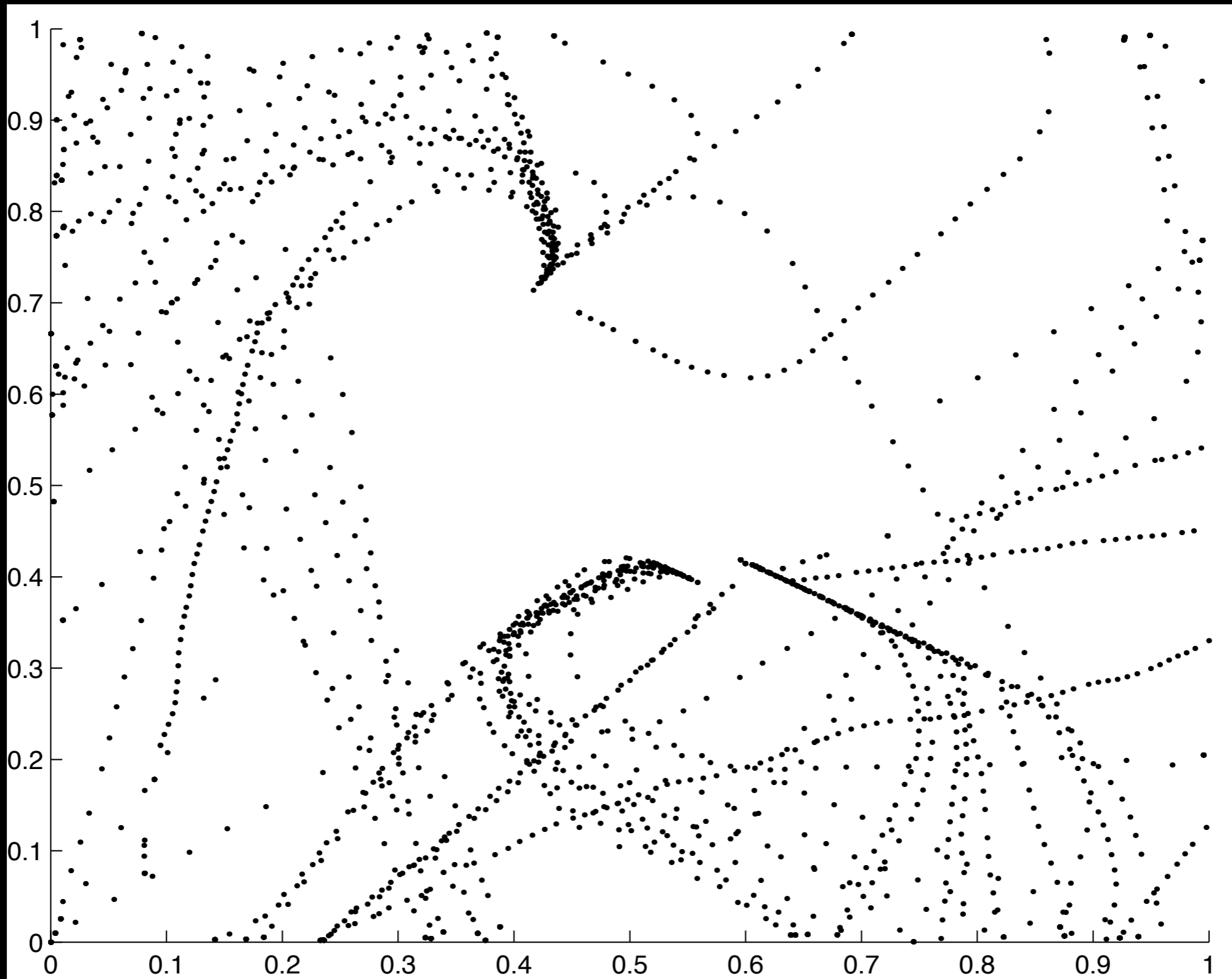
**Efficient inference for distributions on permutations [Huang, Guestrin, Guibas, NIPS 2007]**

**The skew spectrum of graphs [K., Borgwardt, ICML 2008]**

**Characteristic kernels on groups and semigroups [Fukumizu, Sriperumbudur, Gretton, Scholkopf, NIPS 2008]**

**Inferring rankings under constrained sensing [Jagabathula, Shah, NIPS 2008]**

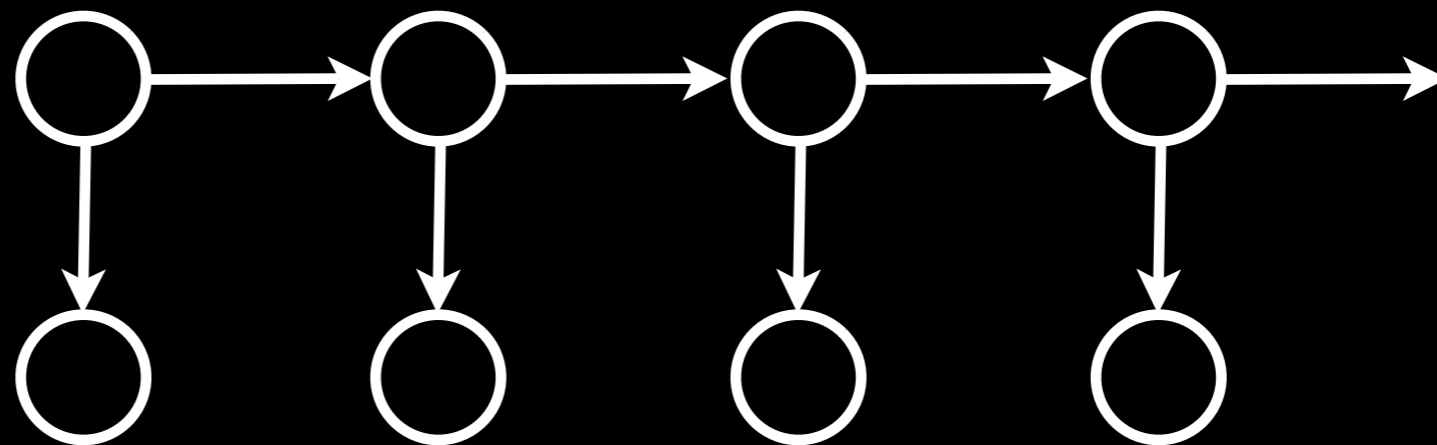
# I. Band-limited approximations



<http://www4.passur.com/jfk.html>

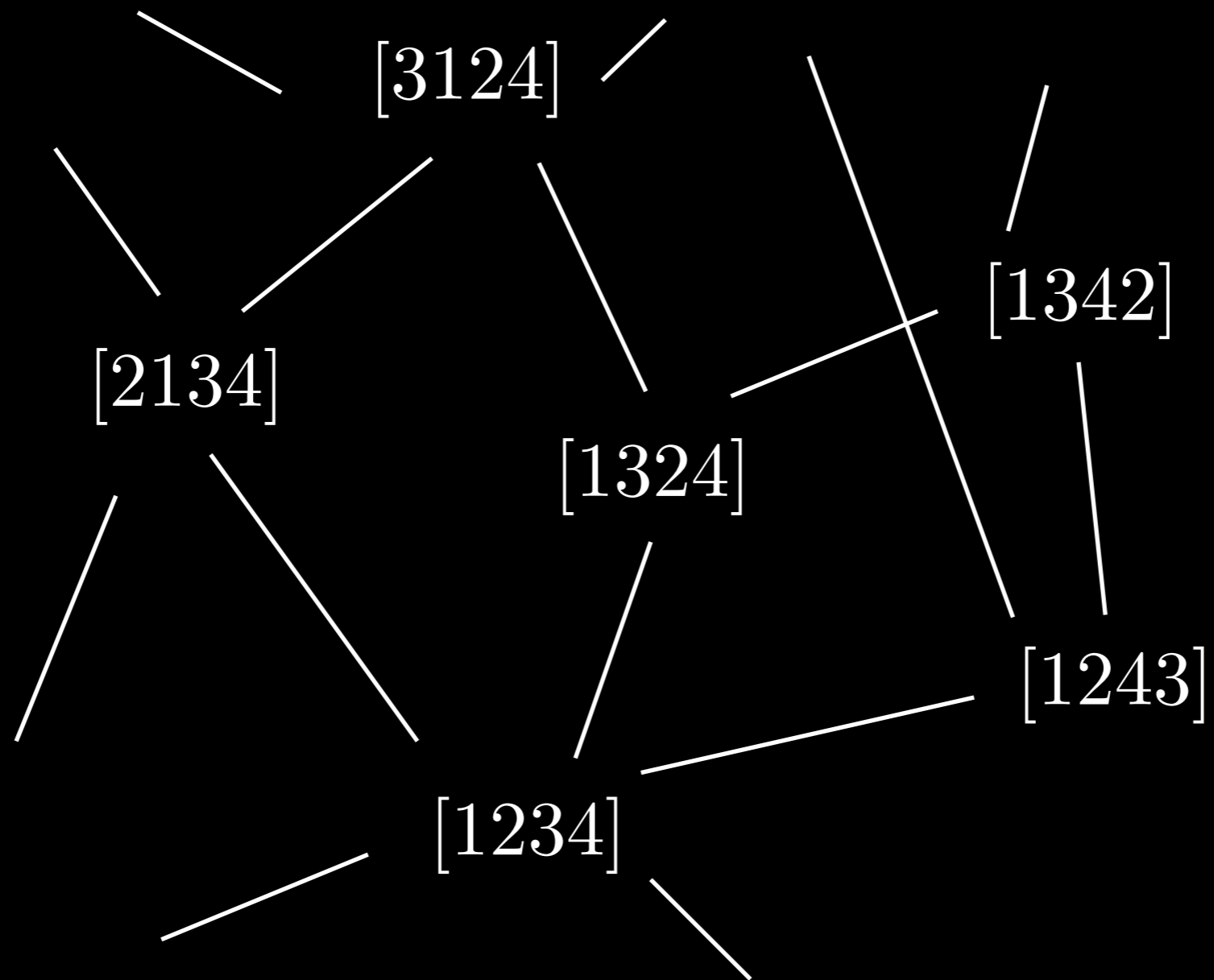
$n$	$n!$
2	2
3	6
4	24
5	120
6	720
7	5040
8	40320
9	362880
10	3628800
11	39916800
12	$4.8 \cdot 10^8$
⋮	⋮
15	$1.3 \cdot 10^{12}$
⋮	⋮
20	$2.4 \cdot 10^{18}$

$$p_t(x_1, x_2, \dots, x_n, \sigma)$$



What does it mean for a function over permutations to be smooth?

Cayley graph:



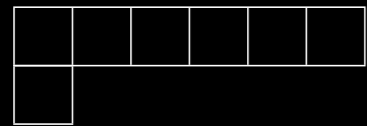
## Theorem

The eigenvectors of the Laplacian of the Cayley graph are the vectors  $v_{\lambda,i,j}(\sigma) = [\rho_{\lambda}(\sigma)]_{i,j}$  and the corresponding eigenvalues are

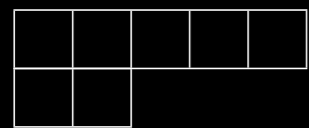
$$\alpha_{\lambda} = \binom{n}{2} \left( 1 - \frac{\text{tr} [\rho_{\lambda}((1, 2))]}{d_{\lambda}} \right)$$



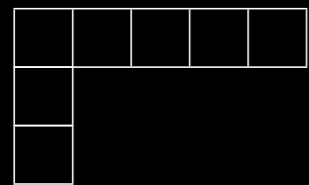
$$d = 1$$



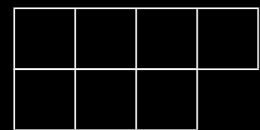
$$d = n - 1$$



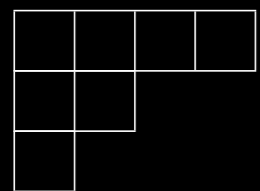
$$d = n(n - 3)/2$$



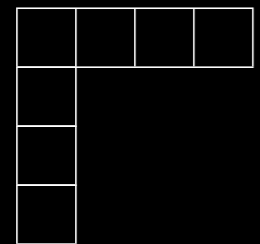
$$d = (n - 1)(n - 2)/2$$



$$d = n(n - 1)(n - 5)/6$$



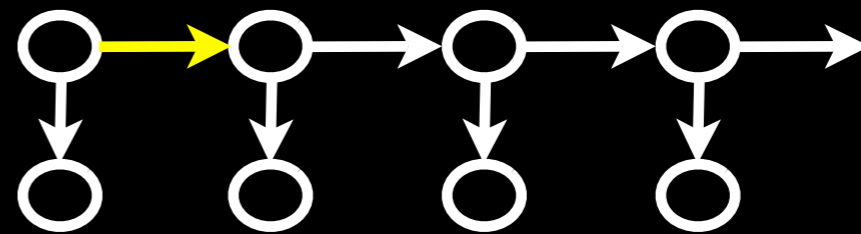
$$d = n(n - 2)(n - 4)/3$$



$$d = (n - 1)(n - 2)(n - 3)/6$$



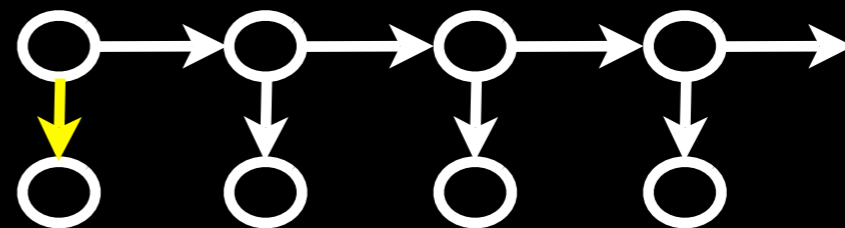
## 1. Noise/filtering/prediction



$$\hat{f}_{t+1}(\rho_\lambda) = \hat{g}(\rho_\lambda) \cdot \hat{f}_t(\rho_\lambda).$$

$$O(D^2)$$

## 2. Observations/conditioning



$$p(O_{a \rightarrow i}^\pi | \sigma) = \begin{cases} \pi & \text{if } \sigma(a) = i, \\ (1 - \pi)/(n - 1) & \text{if } \sigma(a) \neq i. \end{cases}$$

$$f_{t+1}(\sigma) = \frac{p(O_{a \rightarrow i} | \sigma) f_t(\sigma)}{\sum_{\sigma' \in \Sigma} p(O_{a \rightarrow i} | \sigma') f(\sigma')} \cdot \quad O(D^2 n)$$

## 3. Inference

$$\mathbb{P}\{\sigma(n) = i\} = \sum_{\sigma(a)=i} f(\sigma) = \hat{f}_{i,n}(\rho_{(n-1)}) \quad O(n^3)$$

[K., Howard, Jebara, 2007]

- Enforcing positivity constraints by projection.
- More general, but somewhat more expensive updates requiring Clebsch-Gordan decomposition.

---

## Efficient Inference for Distributions on Permutations

---

Jonathan Huang  
Carnegie Mellon University  
jchi@cs.cmu.edu

Carlos Guestrin  
Carnegie Mellon University  
guestrin@cs.cmu.edu

Leonidas Guibas  
Stanford University  
guibas@cs.stanford.edu

### Abstract

Permutations are ubiquitous in many real world problems, such as voting, rankings and data association. Representing uncertainty over permutations is challenging, since there are  $n!$  possibilities, and typical compact representations such as graphical models cannot efficiently capture the mutual exclusivity constraints associated with permutations. In this paper, we use the “low-frequency” terms of a Fourier decomposition to represent such distributions compactly. We present *Kronecker conditioning*, a general and efficient approach for maintaining these distributions directly in the Fourier domain. Low order Fourier-based approximations can lead to functions that do not correspond to valid distributions. To address this problem, we present an efficient quadratic program defined directly in the Fourier domain to project the approximation onto a relaxed form of the marginal polytope. We demonstrate the effectiveness of our approach on a real camera-based multi-people tracking setting.

### 1 Introduction

Permutations arise naturally in a variety of real situations such as card games, data association problems, ranking analysis, etc. As an example, consider a sensor network that tracks the positions of  $n$  people, but can only gather identity information when they walk near certain sensors. Such mixed-modality sensor networks are an attractive alternative to exclusively using sensors which can measure identity because they are potentially cheaper, easier to deploy, and less intrusive. See [1] for a real deployment. A typical tracking system maintains tracks of  $n$  people and the identity of the person corresponding to each track. What makes the problem difficult is that identities can be confused when tracks cross in what we call mixing events. Maintaining accurate track-to-identity assignments in the face of these ambiguities based on identity measurements is known as the *Identity Management Problem* [2], and is known to be *NP-hard*. Permutations pose a challenge for probabilistic inference, because distributions on the group of permutations on  $n$  elements require storing at least  $n! - 1$  numbers, which quickly becomes infeasible as  $n$  increases. Furthermore, typical compact representations, such as graphical models, cannot capture the mutual exclusivity constraints associated with permutations.

Diaconis [3] proposes maintaining a small subset of Fourier coefficients of the actual distribution allowing for a principled tradeoff between accuracy and complexity. Schumitsch et al. [4] use similar ideas to maintain a particular subset of Fourier coefficients of the log probability distribution. Kondor et al. [5] allow for general sets of coefficients, but assume a restrictive form of the observation model in order to exploit an efficient FFT factorization. The main contributions of this paper are:

- A new, simple and general algorithm, *Kronecker Conditioning*, which performs all probabilistic inference operations completely in the Fourier domain. Our approach is general, in the sense that it can address any transition model or likelihood function that can be represented in the Fourier domain, such as those used in previous work, and can represent the probability distribution with any desired set of Fourier coefficients.
- We show that approximate conditioning can sometimes yield Fourier coefficients which do not correspond to any valid distribution, and present a method for projecting the result back onto a relaxation of the marginal polytope.
- We demonstrate the effectiveness of our approach on a real camera-based multi-people tracking setting.

[Huang, Guestrin, Guibas, 2007]

## 2. Invariances

Recall:  $\widehat{f^z}(\rho) = \rho(z)\widehat{f}(\rho)$        $f^z(x) = f(z^{-1}x)$

Power spectrum:  $\widehat{a}(\rho) = \widehat{f}(\rho)^\dagger \widehat{f}(\rho)$

Bispectrum:

$$b(\rho_1, \rho_2) = C^\dagger (\widehat{f}(\rho_1) \otimes \widehat{f}(\rho_2))^\dagger C \bigoplus_{\rho} \widehat{f}(\rho)$$

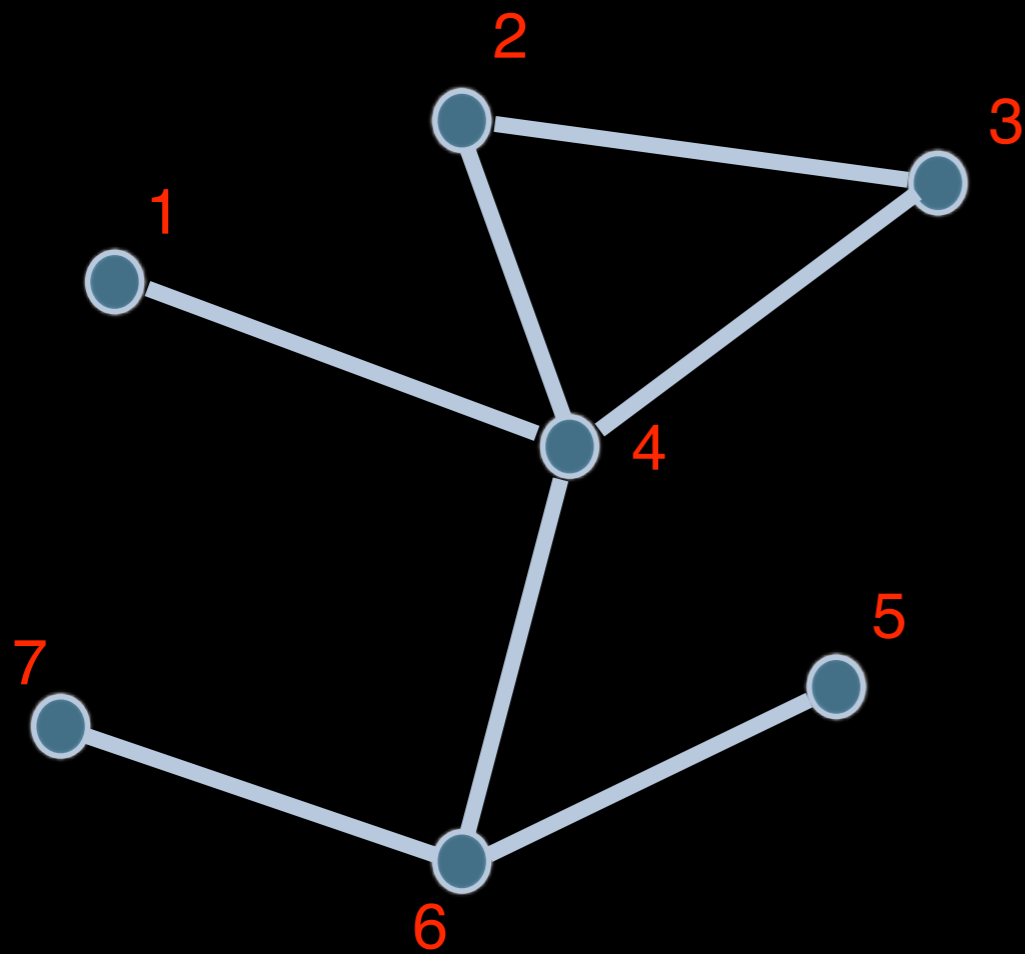
$$\rho_1(z) \otimes \rho_2(z) = C \left[ \bigoplus_{\rho} \rho(z) \right] C^\dagger$$

Skew spectrum:

$$\hat{q}_z(\rho) = \hat{r}_z(\rho)^\dagger \cdot \hat{f}(\rho) \quad r_z(x) = f(xz)f(x) \quad z \in G$$

(unitarily equivalent, but easier to compute)

[K. 2008]



$$A = \begin{pmatrix} 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

$$[A^\pi]_{\pi(i),\pi(j)} = [A]_{i,j}$$

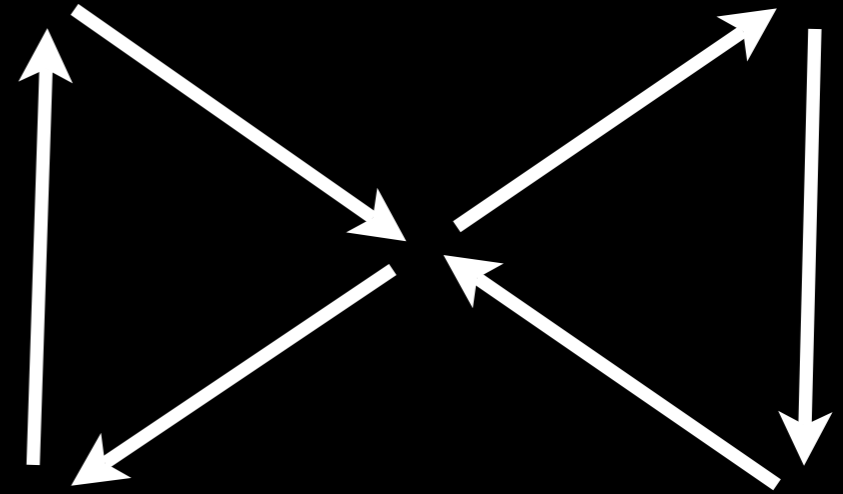
[K., Borgwardt, 2008]

$$\hat{f}(\rho_1) = \boxed{(6)}$$

$$\hat{f}(\rho_2) = \begin{pmatrix} \boxed{-0.25} & \boxed{-0.323} & 0 & 0 \\ \boxed{-0.323} & \boxed{-0.417} & 0 & 0 \\ \boxed{0.913} & \boxed{1.18} & 0 & 0 \\ \boxed{0} & \boxed{0} & 0 & 0 \end{pmatrix}$$

$$\hat{f}(\rho_3) = \begin{pmatrix} \boxed{1.33} & 0 & 0 & 0 & 0 \\ \boxed{0.471} & 0 & 0 & 0 & 0 \\ \boxed{0} & 0 & 0 & 0 & 0 \\ \boxed{0.816} & 0 & 0 & 0 & 0 \\ \boxed{0} & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\hat{f}(\rho_4) = \begin{pmatrix} \boxed{-1.67} & 0 & 0 & 0 & 0 & 0 \\ \boxed{1.18} & 0 & 0 & 0 & 0 & 0 \\ \boxed{0} & 0 & 0 & 0 & 0 & 0 \\ \boxed{-0.913} & 0 & 0 & 0 & 0 & 0 \\ \boxed{0} & 0 & 0 & 0 & 0 & 0 \\ \boxed{-2.24} & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$



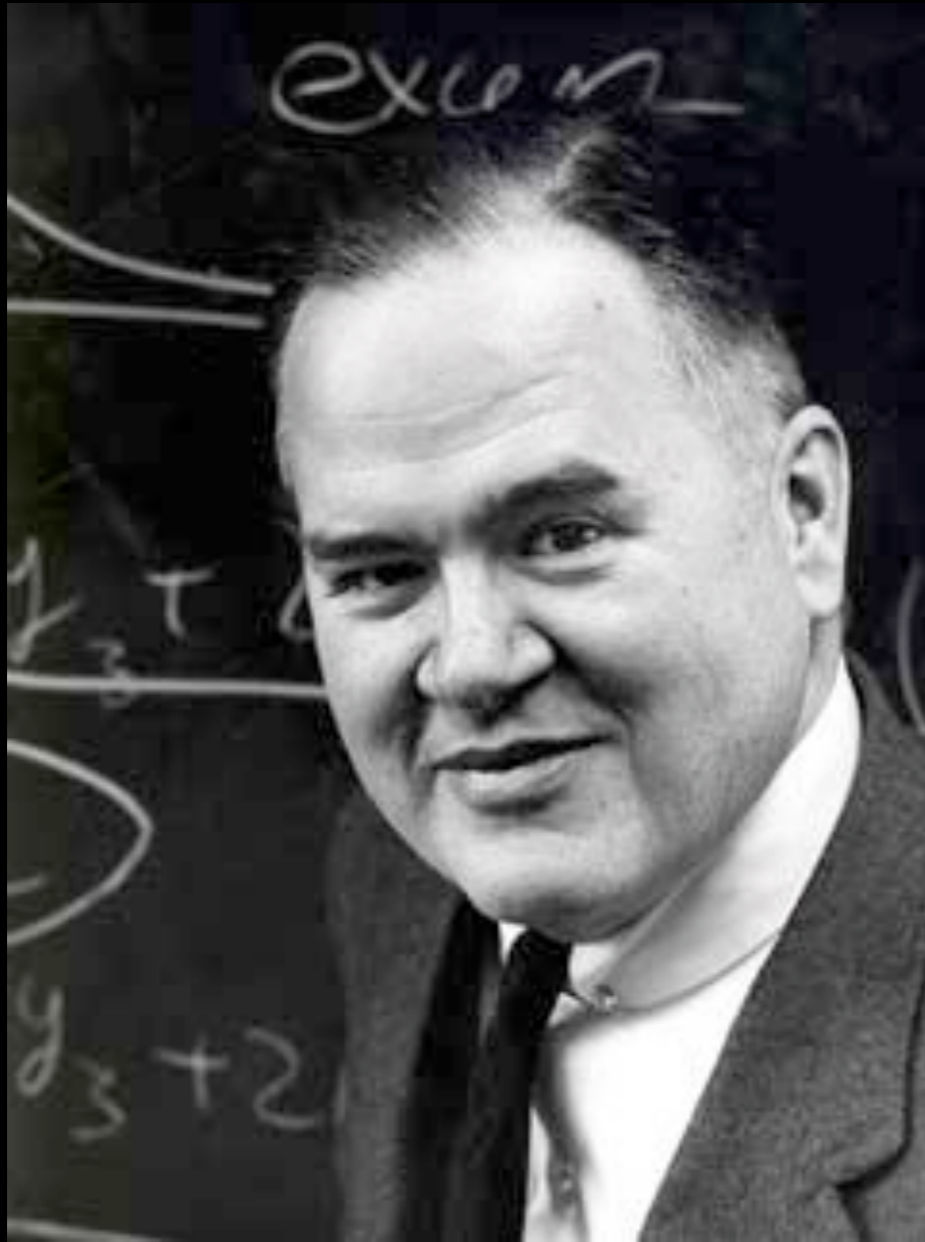
## Proposition [K. & Borgwardt 2008]

The reduced skew spectrum of a graph contains 49 scalar components and may be computed in  $O(n^3)$  operations.

	MUTAG	ENZYME	NCI1	NCI109
Number of instances/classes	600/6	188/2	4110/2	4127/2
Max. number of nodes	28	126	111	111
Reduced skew spectrum	<b>88.61</b> (0.21)	25.83 (0.34)	<b>62.72</b> (0.05)	<b>62.62</b> (0.03)
Random walk kernel	71.89 (0.66)	14.97 (0.28)	51.30 (0.23)	53.11 (0.11)
Shortest path kernel	81.28 (0.45)	<b>27.53</b> (0.29)	61.66 (0.10)	62.35 (0.13)



# 3. Algorithms



"An algorithm for the machine calculation of complex Fourier series," James W. Cooley and John W. Tukey, *Math. Comput.* **19**, 297–301 (1965)

John Wilder Tukey  
1915-2000

“Fast generalized Fourier transforms” [Clausen, 1989]

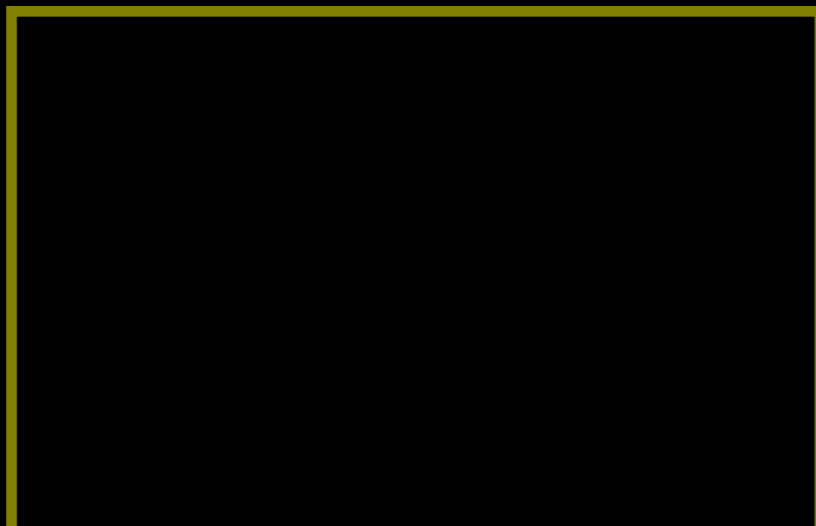
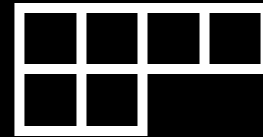
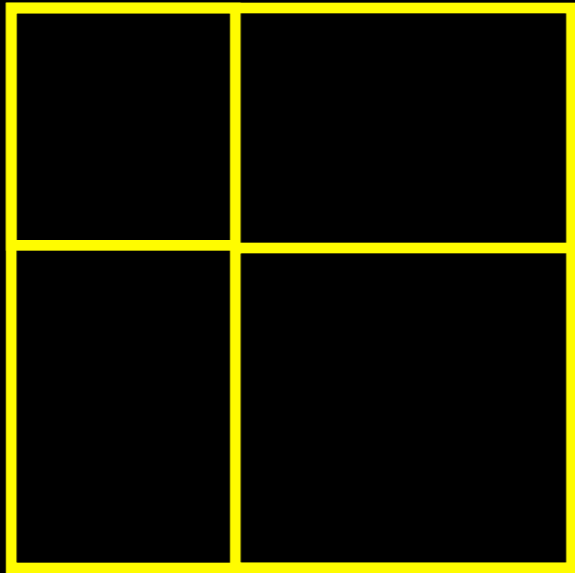
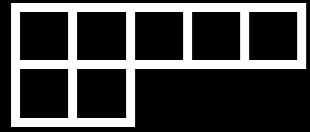
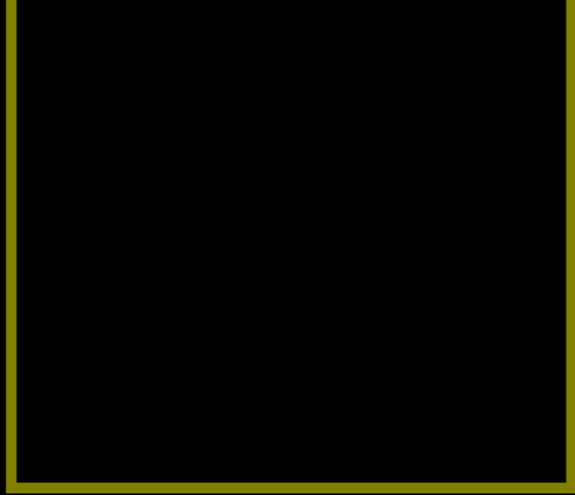
FFT for wreath product groups and many others by Maslen, Rockmore, Healy,...



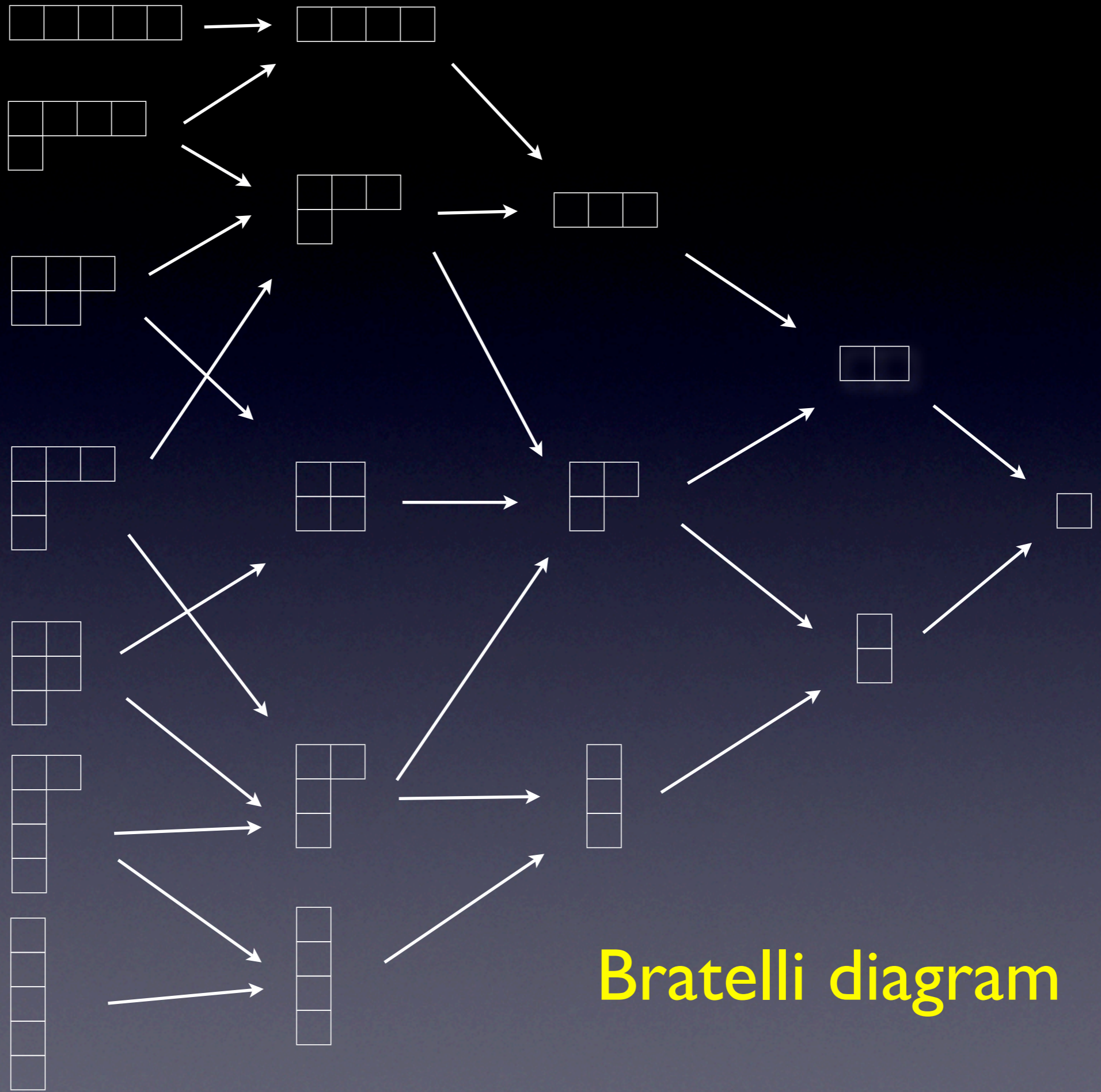
Michael Clausen



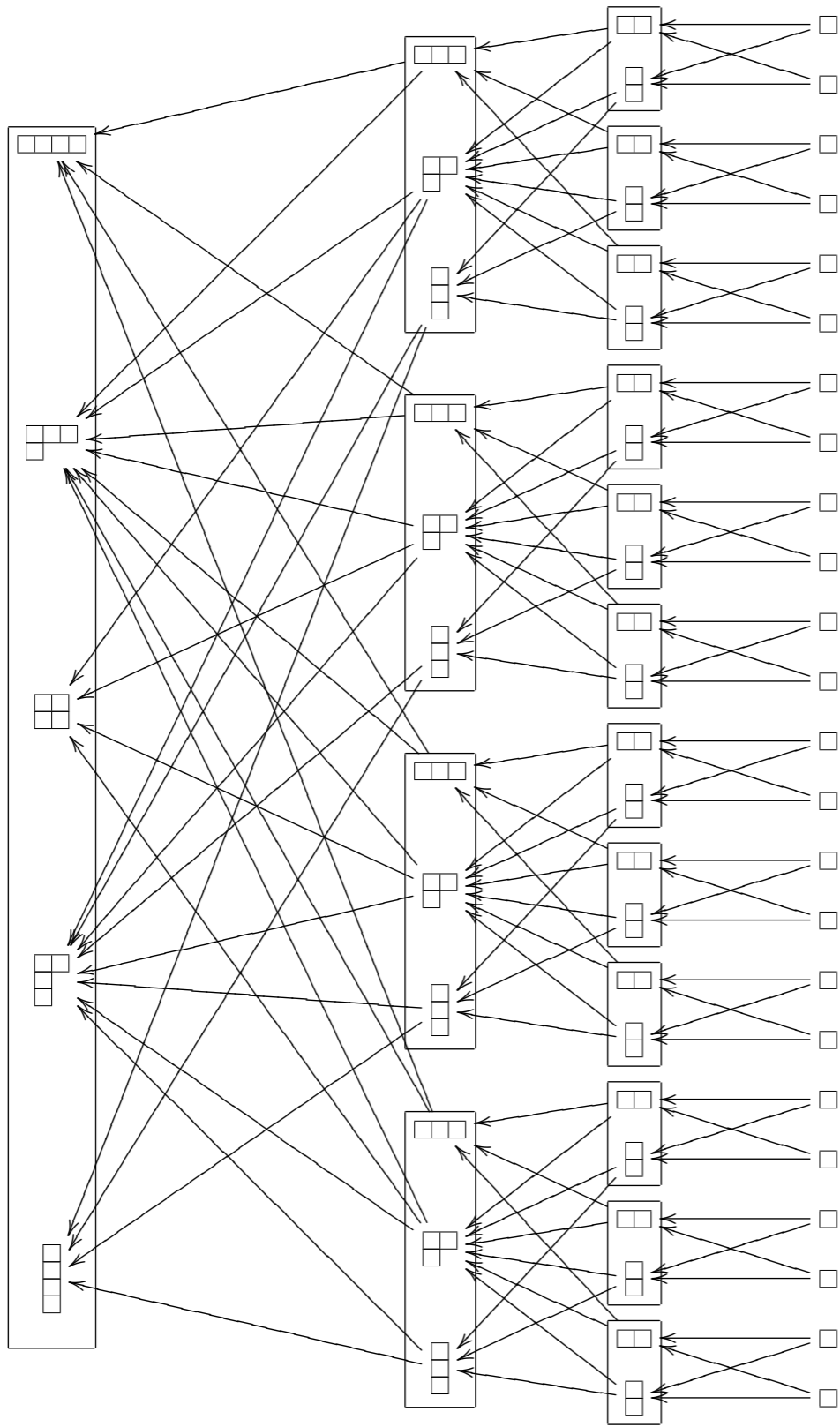






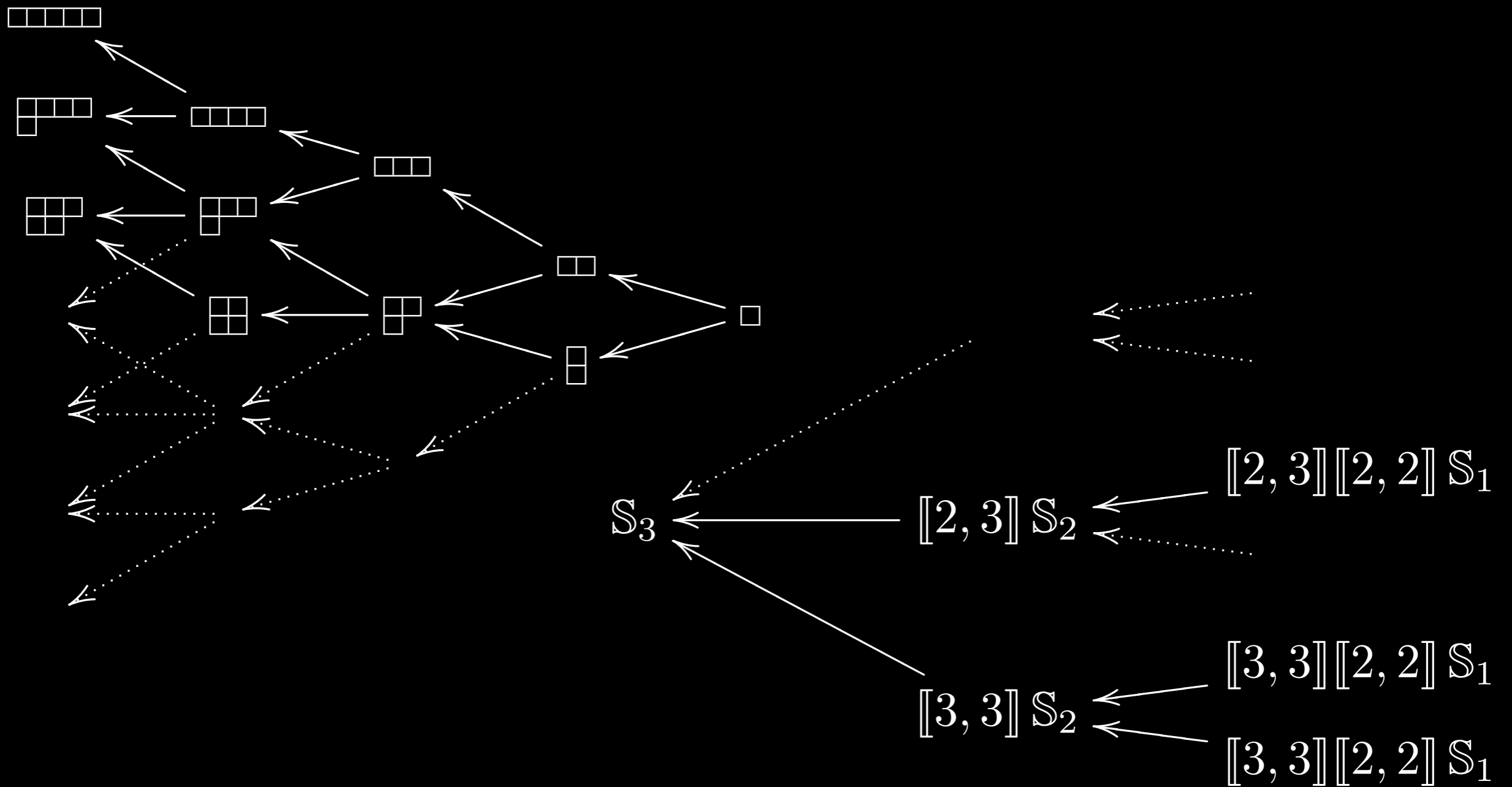


**Bratelli diagram**



$$n!^2 \rightarrow \frac{(n+1)n(n-1)}{3} n!$$

# Sparse transforms





# $S_n$ ob

A C++ library for fast Fourier transforms on the symmetric group.

author: Risi Kondor, Columbia University ([risi@cs.columbia.edu](mailto:risi@cs.columbia.edu))

Development version as of August 23, 2006 (unstable!):

Documentation: [\[ps\]](#)[\[pdf\]](#)

C++ source code: [\[directory\]](#)

BiBTeX entry: [\[bib\]](#)

Entire package: [\[tar.gz\]](#)

ALL SOFTWARE ON THIS PAGE IS DISTRIBUTED UNDER THE TERMS OF THE GNU  
GENERAL PUBLIC LICENSE [\[site\]](#)

## References:

1. Michael Clausen: **Fast generalized Fourier transforms**. Theoretical Computer Science **67(1)**: 55-63, 1989.
2. David K. Maslen and Daniel N. Rockmore: **Generalized FFTs --- a survey of some recent results**. Proceedings of the DIMACS Workshop on Groups and Computation, 1997. [\[ps\]](#)
3. K.-J. Kueh, T. Olson, D. Rockmore and K.-S. Tan: **Nonlinear approximation theory on finite**

<http://www.gatsby.ucl.ac.uk/~risi/SnOB>

# Conclusions

Harmonic analysis generalizes naturally to non-commutative groups of transformations.

Strong connections to well developed branches of mathematics.

Rich algorithmic side: FFTs, etc.