# Approaches to bounding the exponent of matrix multiplication
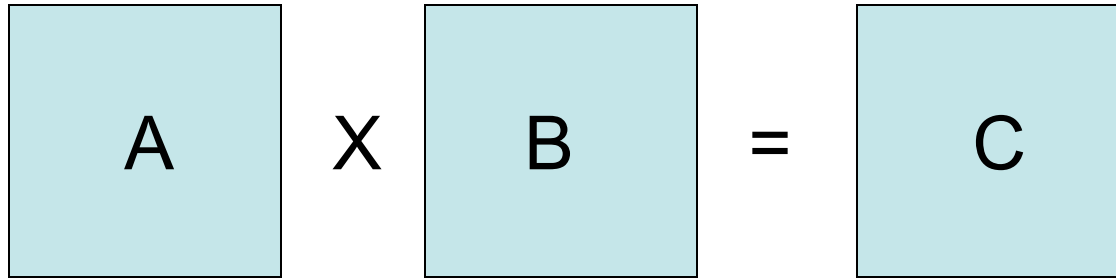
## Chris Umans

Caltech

Based on joint work with Noga Alon, Henry Cohn, Bobby Kleinberg, Amir Shpilka, Balazs Szegedy
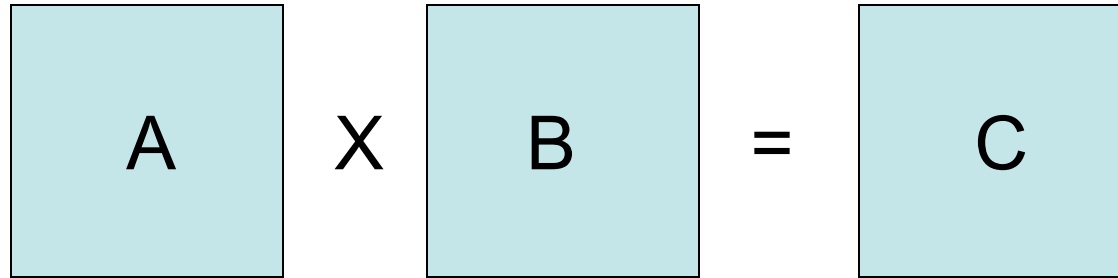
Simons Institute Sept. 17, 2014

# Introduction

A   X   B   =   C

- Standard method: $O(n^3)$ operations
- Strassen (1969): $O(n^{2.81})$ operations

# Introduction

$$A \times B = C$$

- Standard method: $O(n^3)$ operations
- Strassen (1969): $O(n^{2.81})$ operations

> The exponent of matrix multiplication: smallest number $\omega$ such that for all $\varepsilon > 0$
>
> $O(n^{\omega + \varepsilon})$ operations suffice

# History

- Standard algorithm                                     $\omega \leq 3$
- Strassen (1969)                                        $\omega < 2.81$
- Pan (1978)                                             $\omega < 2.79$
- Bini; Bini et al. (1979)                              $\omega < 2.78$
- Schönhage (1981)                                       $\omega < 2.55$
- Pan; Romani; Coppersmith
   + Winograd (1981-1982)                                $\omega < 2.50$
- Strassen (1987)                                        $\omega < 2.48$
- Coppersmith + Winograd (1987)                          $\omega < 2.375$
- Stothers (2010)                                        $\omega < 2.3737$
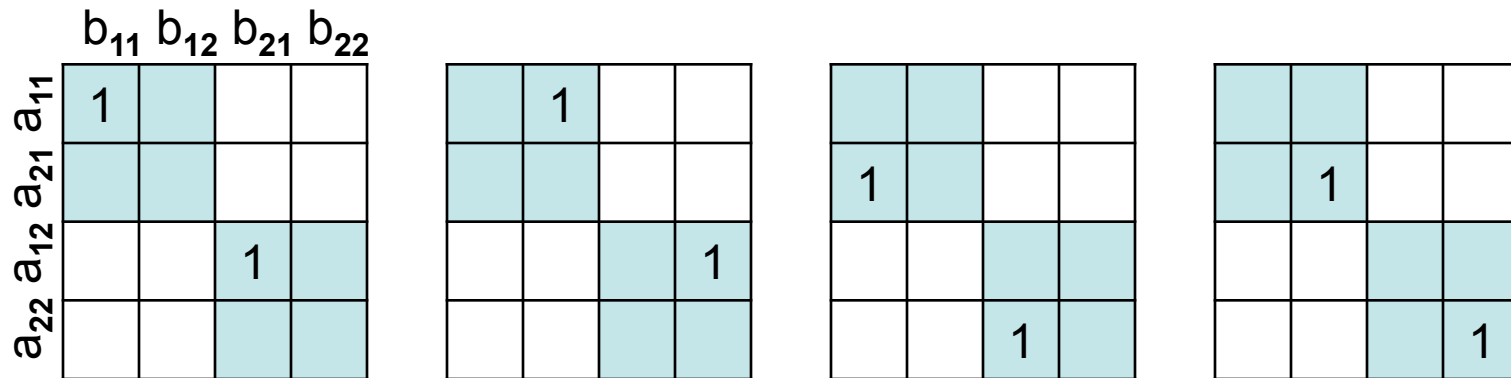- Williams (2011)                                        $\omega < 2.3729$
- Le Gall (2014)                                         $\omega < 2.37286$

# Outline

1. main ideas from Strassen 1969 through Le Gall 2014

2. approach via embedding into semi-simple algebra multiplication

   – groups

   – coherent configurations/association schemes

# The matrix multiplication tensor

<n,n,n> is a $n^2$ x $n^2$ x $n^2$ tensor described by trilinear form $\sum_{i,j,k} X_{i,j} Y_{j,k} Z_{k,i}$
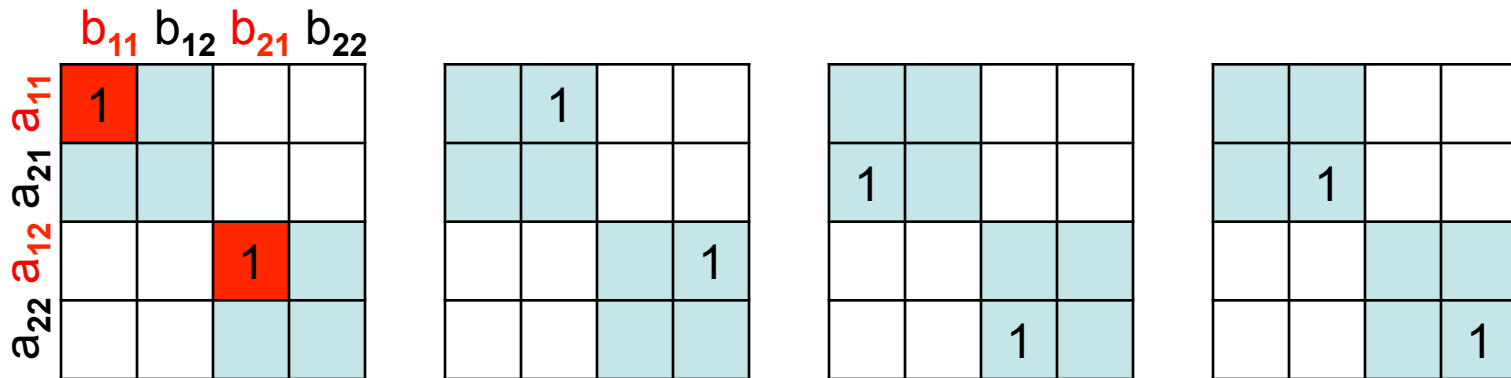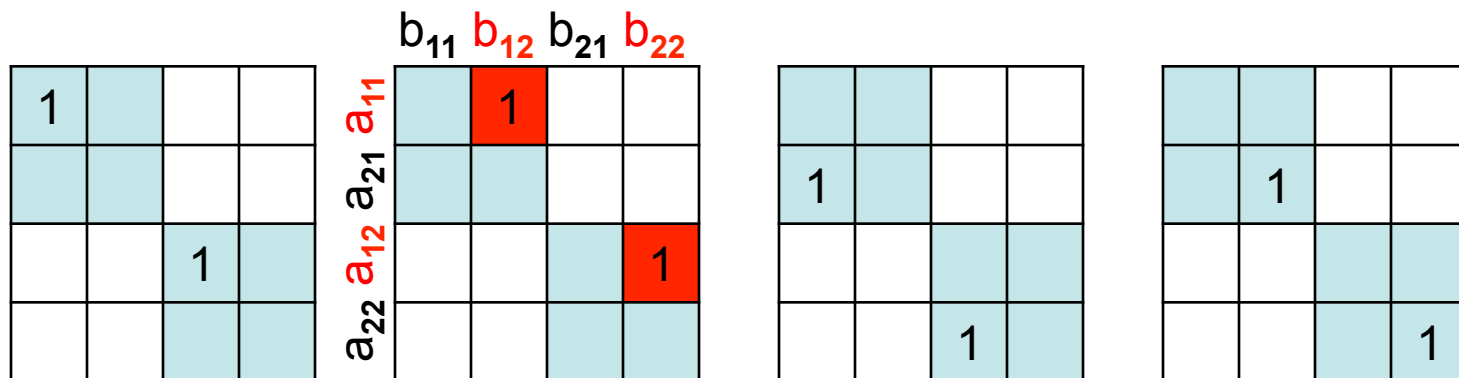
$$\begin{array}{|c|c|} \hline a_{11} & a_{12} \\ \hline a_{21} & a_{22} \\ \hline \end{array} \quad \times \quad \begin{array}{|c|c|} \hline b_{11} & b_{12} \\ \hline b_{21} & b_{22} \\ \hline \end{array} \quad = \quad \begin{array}{|c|c|} \hline c_{11} & c_{12} \\ \hline c_{21} & c_{22} \\ \hline \end{array}$$

|  | $b_{11}$ | $b_{12}$ | $b_{21}$ | $b_{22}$ |
|---|---|---|---|---|
| $a_{11}$ | 1 |  |  |  |
| $a_{21}$ |  |  |  |  |
| $a_{12}$ |  |  | 1 |  |
| $a_{22}$ |  |  |  |  |

|  |  |  |  |
|---|---|---|---|
|  | 1 |  |  |
|  |  |  |  |
|  |  |  | 1 |
|  |  |  |  |

|  |  |  |  |
|---|---|---|---|
|  |  |  |  |
| 1 |  |  |  |
|  |  | 1 |  |
|  |  |  |  |

|  |  |  |  |
|---|---|---|---|
|  |  |  |  |
|  | 1 |  |  |
|  |  |  | 1 |
|  |  |  |  |

# The matrix multiplication tensor

<n,n,n> is a $n^2$ x $n^2$ x $n^2$ tensor described by trilinear form $\sum_{i,j,k} X_{i,j} Y_{j,k} Z_{k,i}$

| $a_{11}$ | $a_{12}$ |
|----------|----------|
| $a_{21}$ | $a_{22}$ |

x

| $b_{11}$ | $b_{12}$ |
|----------|----------|
| $b_{21}$ | $b_{22}$ |

=

| $c_{11}$ | $c_{12}$ |
|----------|----------|
| $c_{21}$ | $c_{22}$ |

# The matrix multiplication tensor

$\langle$n,n,n$\rangle$ is a $n^2$ x $n^2$ x $n^2$ tensor described by trilinear form $\sum_{i,j,k} X_{i,j} Y_{j,k} Z_{k,i}$
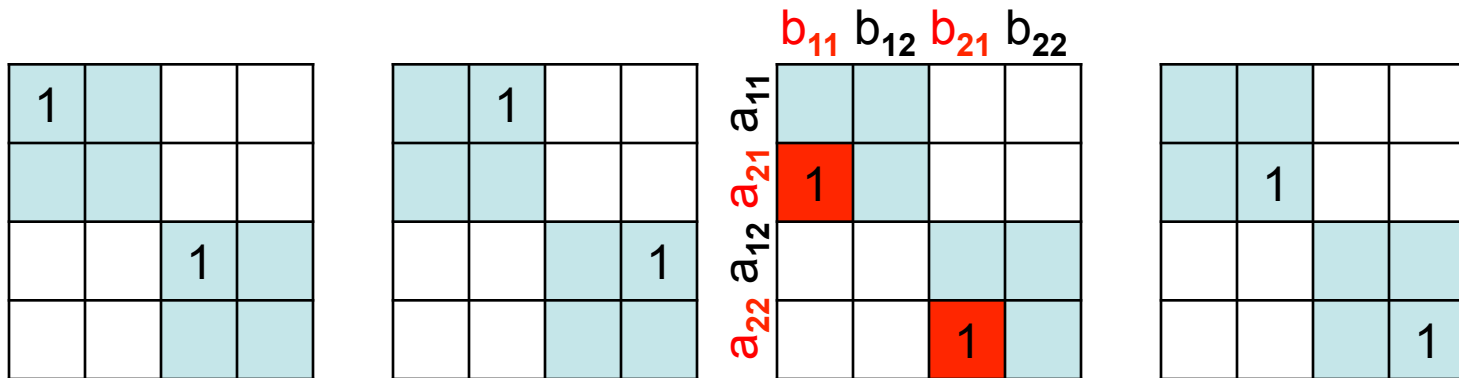
# The matrix multiplication tensor

$\langle n,n,n \rangle$ is a $n^2$ x $n^2$ x $n^2$ tensor described by trilinear form $\sum_{i,j,k} X_{i,j} Y_{j,k} Z_{k,i}$

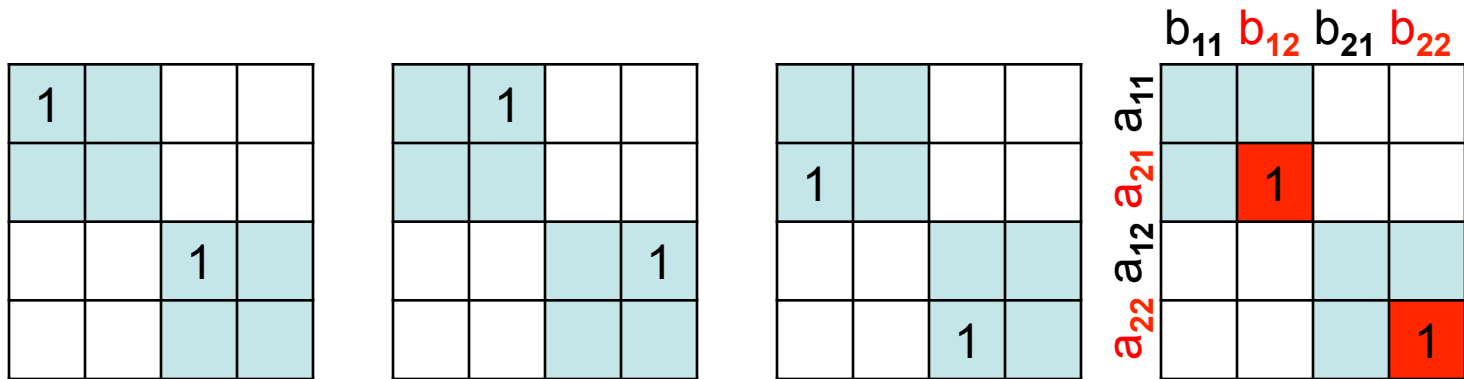# The matrix multiplication tensor

<n,n,n> is a $n^2$ x $n^2$ x $n^2$ tensor described by trilinear form $\sum_{i,j,k} X_{i,j} Y_{j,k} Z_{k,i}$

# The matrix multiplication tensor

<n,m,p> is a nm £ mp £ pn tensor
described by trilinear form $\sum_{i,j,k} X_{i,j} Y_{j,k} Z_{k,i}$

$$m$$

$$n \quad A \quad \times \quad m \quad B \quad = \quad n \quad C$$

$$p \qquad p$$

Each of
np slices of
<n,m,p>:

1

À m !

...

1

# Strategies
# for upper bounding the rank
# of the
# matrix multiplication tensor

# Upper bounds on rank

- Observation: $\langle n,n,n\rangle^{-i} = \langle n^i, n^i, n^i\rangle$

  $)\ R(\langle n^i, n^i, n^i\rangle) \cdot R(\langle n,n,n\rangle)^i$

- Strategy I: bound rank for small n by hand
  - $R(\langle 2,2,2\rangle) = 7$          ! < 2.81
  - $R(\langle 3,3,3\rangle)$ 2 [19..23]        (worse bound)

  - even computer search infeasible…

# Upper bounds on rank

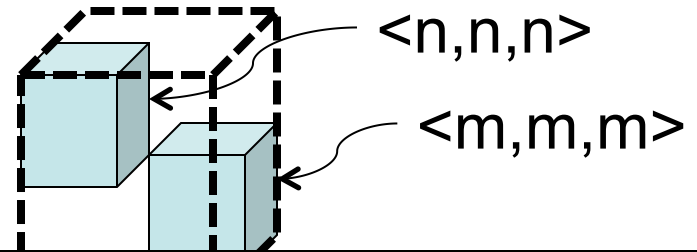- Border rank = rank of sequence of tensors approaching target tensor entrywise



rank = 3
border rank = 2:

- Strategy II: bound *border rank* for small n

- Lemma: $\underline{R}(<n,n,n>) < r ) !< \log_n r$
  - $\underline{R}(<2,2,3>) \cdot 10$        $!< 2.79$

# Upper bounds on rank

- Direct sum of tensors
  <n,n,n> © <m,m,m>

  (multiple matrix multi



<n,n,n>

<m,m,m>

"Asymptotic Sum Inequality" and example (Schönhage 1981)

- Strategy III: bound (border) rank of *direct sums* of small matrix multiplication tensors

$\underline{R}(<n_1,n_1,n_1> © … © <n_k,n_k,n_k>) < r ) \sum_i n_i^! < r$

  – $\underline{R}(<4,1,3> © <1,6,1>) \cdot 13$     ! < 2.55

# Upper bounds on rank

- ## Strategy IV: Strassen "laser method"
  - tensor with "coarse structure" of MM and "fine structure" components isomorphic to MM

(many independent MMs in high tensor powers)



coarse structure
<1,2,1>

fine = scalar **x** row vector
col vector **x** scalar

# Upper bounds on rank

- **Strategy IV**: Strassen "laser method"
  - tensor with "coarse structure" of MM and "fine structure" components isomorphic to MM

(many independent MMs in high tensor powers)
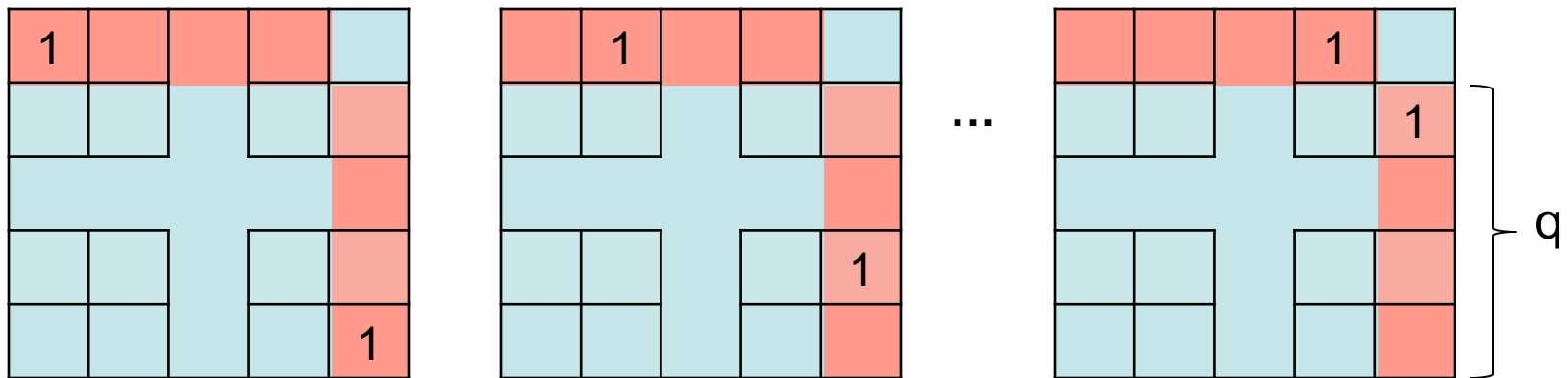


border rank = $q + 1$;          $q = 5$ yields **!** < 2.48

# Upper bounds on rank

- Coppersmith-Winograd and beyond: border rank of this tensor is q+2:

$$\sum_{i=1\ldots q} X_0 Y_i Z_i + X_i Y_0 Z_i + X_i Y_i Z_0 +$$

$$X_0 Y_0 Z_{q+1} + X_0 Y_{q+1} Z_0 + X_{q+1} Y_0 Z_0$$

- 6 "pieces": target proportions in high tensor power affect # and size of independent MMs
- q = 6 yields ! < 2.388

# Upper bounds on rank

- **Coppersmith-Winograd and beyond:** analyze tensor powers of this tensor

$$T_q = \sum_{i=1\ldots q} X_0 Y_i Z_i + X_i Y_0 Z_i + X_i Y_i Z_0 +$$
$$X_0 Y_0 Z_{q+1} + X_0 Y_{q+1} Z_0 + X_{q+1} Y_0 Z_0$$

| Tensor power | # "pieces" | bound | reference |
| --- | --- | --- | --- |
| 2 | 36 | 2.375 | C-W |
| 4 | 1296 | 2.3737 | Stothers |
| 8 | 1679616 | 2.3729 | Williams |
| 16 | 2.82 x 10^12 | 2.3728640 | Le Gall |
| 32 | 7.95 x 10^24 | 2.3728639 | Le Gall |

# Upper bounds on rank

- ## Coppersmith-Winograd and beyond

| Tensor power | # pieces | bound | reference |
|---|---|---|---|
| 2 | 36 | 2.375 | C-W |
| 4 | 1296 | 2.3737 | Stothers |
| 8 | 1679616 | 2.3729 | Williams |
| 16 | 2.82 x 10^12 | 2.3728640 | Le Gall |
| 32 | 7.95 x 10^24 | 2.3728639 | Le Gall |

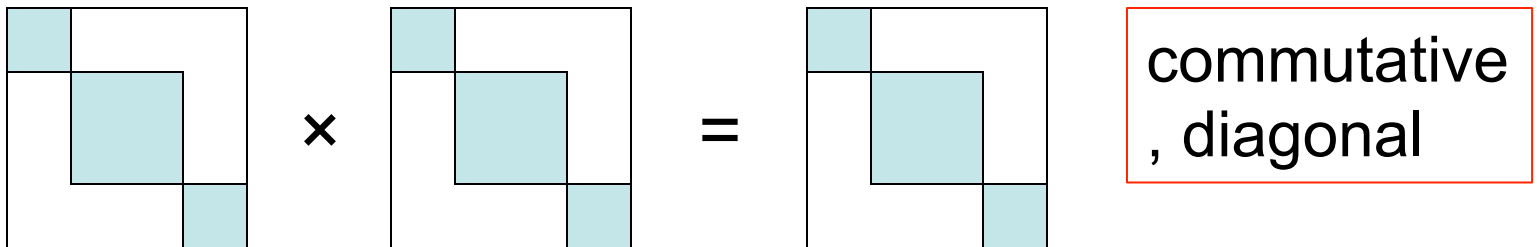- ## Ambainis-Filmus 2014: N-th tensor power cannot beat bound of 2.3078

# A different approach

- ## So far...
  - bound border rank of small tensor (by hand)
  - asymptotic bound from high tensor powers

- ## Disadvantages
  - limited universe of "starting" tensors
  - high tensor powers hard to analyze

matrix multiplication

via groups and

coherent configurations /

association schemes

# The general approach

- Cohn-Umans 2003, 2012:
    - *embed* n x n matrix multiplication into semi-simple algebra multiplication
    - semi-simple: isomorphic to block-diagonal MM



commutative, diagonal

    - key hope: "nice basis" w/ combinatorial structure
    - reduce n x n MM to smaller MMs; recurse

# The Group Algebra

- given finite group G, <span style="color:red">group algebra</span> C[G] has elements $\Sigma_g\, a_g g$

  with multiplication

$$(\Sigma_g a_g g)(\Sigma_h b_h h) = \Sigma_f\, (\Sigma_{gh\,=\,f}\, a_g b_h)f$$

- structure: C[G] $'$ $(C^{d_1 \times d_1}) \times \ldots \times (C^{d_k \times d_k})$
- group elements are "nice basis"

# "Nice basis" embedding:

Subgroups X, Y, Z of G satisfy the
**triple product property**
if for all $x \in X$, $y \in Y$, $z \in Z$ :

$xyz = 1$     iff    $x = y = z = 1$.

# The embedding:

Subsets X, Y, Z of G satisfy the
**triple product property**
if for all $x \in Q(X)$, $y \in Q(Y)$, $z \in Q(Z)$:

$xyz = 1$      iff    $x = y = z = 1$.

$\underline{\mathbf{A}} = \Sigma a_{x,y} (x\ y^{-1})$        $\underline{\mathbf{B}} = \Sigma b_{y,z} (y\ z^{-1})$

**Claim:** $(AB)_{x,z}$ = coeff. on $(x\ z^{-1})$ in $\underline{\mathbf{A}}*\underline{\mathbf{B}}$.

# The embedding:

Subsets X, Y, Z of G satisfy the
**triple product property**
if for all $x \in Q(X)$, $y \in Q(Y)$, $z \in Q(Z)$:

$xyz = 1$    iff   $x = y = z = 1.$

$\underline{\mathbf{A}} = \Sigma a_{x_1,y_1}(x_1 y_1^{-1})$         $\underline{\mathbf{B}} = \Sigma b_{y_2,z_2}(y_2 z_2^{-1})$
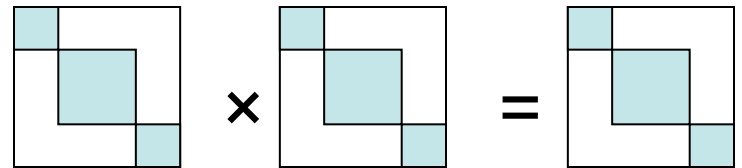
**Claim:** $(AB)_{x_3,z_3}$ = coeff. on $(x_3 z_3^{-1})$ in $\underline{\mathbf{A}}*\underline{\mathbf{B}}$.

$(x_1 y_1^{-1})(y_2 z_2^{-1}) = x_3 z_3^{-1}$   **)** $x_3^{-1} x_1 y_1^{-1} y_2 z_2^{-1} z_3 = 1$

# How many multiplications?

Embedding + structure of C[G] yields bound on rank (´ # multiplications):

- we use m ≤ $\Sigma d_i^3$ mults
- really m = $\Sigma d_i^!$ mults
- *at least* m ≥ $\Sigma d_i^2$ = |G| mults

**First Challenge**: embed k × k matrix multiplication in group of size ¼ $k^2$

# The embedding

**First Challenge**: embed $k \times k$ matrix multiplication in group of size $\frac{1}{4} k^2$

- simple pigeonhole argument:
  - embedding in an abelian group requires group to have size $k^3$

# The triangle construction

**Theorem**: can embed k × k matrix multiplication in <span style="color:red">symmetric group</span> of size $k^{2 + o(1)}$

n objects →



- subgroup X
- subgroup Y
- subgroup Z

need X, Y, Z in $S_n$ all with size $\approx |S_n|^{1/2}$

# The triangle construction

– X moves points within rows

– Y moves points within columns

– Z moves points within diagonals

– want: xyz = 1 $\Rightarrow$ x = y = z = 1

# The triangle construction

**Theorem**: can embed k × k matrix multiplication in <span style="color:red">symmetric group</span> of size $k^{2 + o(1)}$
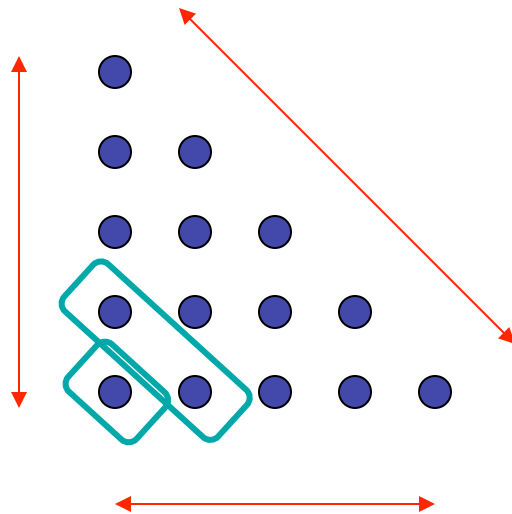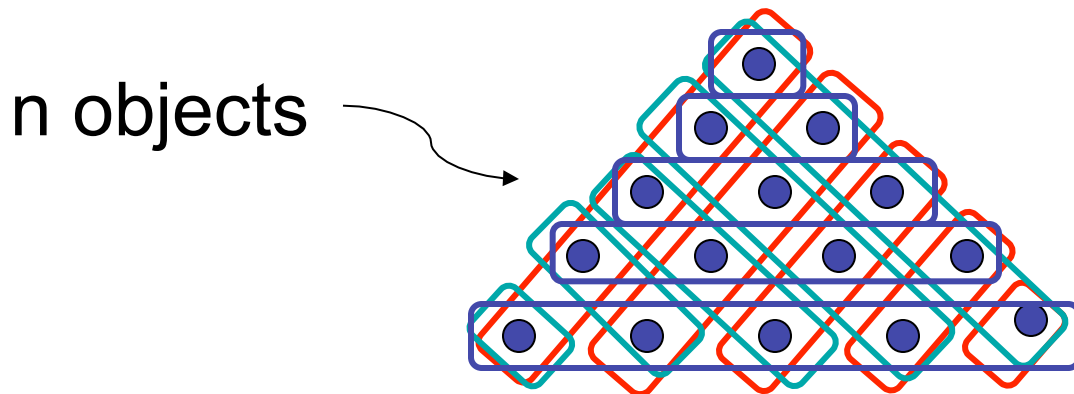
n objects



- subgroup X
- subgroup Y
- subgroup Z

unfortunately, $d_{max} > |X|$ (= $|Y|$ = $|Z|$)

# What should we be aiming for?

**Theorem**: in group G supporting $k \times k$ matrix multiplication with character degrees $d_1, d_2, d_3, \ldots$, we obtain:

$$k^\omega \cdot \sum_i d_i^\omega$$

- If $X, Y, Z \subseteq G$ satisfy T.P.P. and
  - $(|X| \cdot |Y| \cdot |Z|)^{1/3} = k \geq |G|^{1/2 - o(1)}$
  - $d_{max} \cdot |G|^{1/2 - \epsilon}$

  then $\omega = 2$

$$\sum_i d_i^\omega \cdot d_{max}^{\omega - 2} |G|$$

# Constructions in linear groups

- Good candidate family:

    SL(n, q) for fixed dimension n

- In SL(n, R) these three subgroups satisfy the triple product property:

    – upper-triangular with ones on the diagonal

    – lower-triangular with ones on the diagonal

    – the special orthogonal group SO(n, R)

    and dim. of each is ½ dim. of G as n ! 1

# Group algebra approach

- [CKSU 2005] wreath product groups yield :
  - $\omega < 2.48$, $\omega < 2.41$
  - key part of construction is combinatorial
  - two conjectures implying $\omega = 2$

- Main disadvantage:
  - non-trivial results *require* non-abelian groups
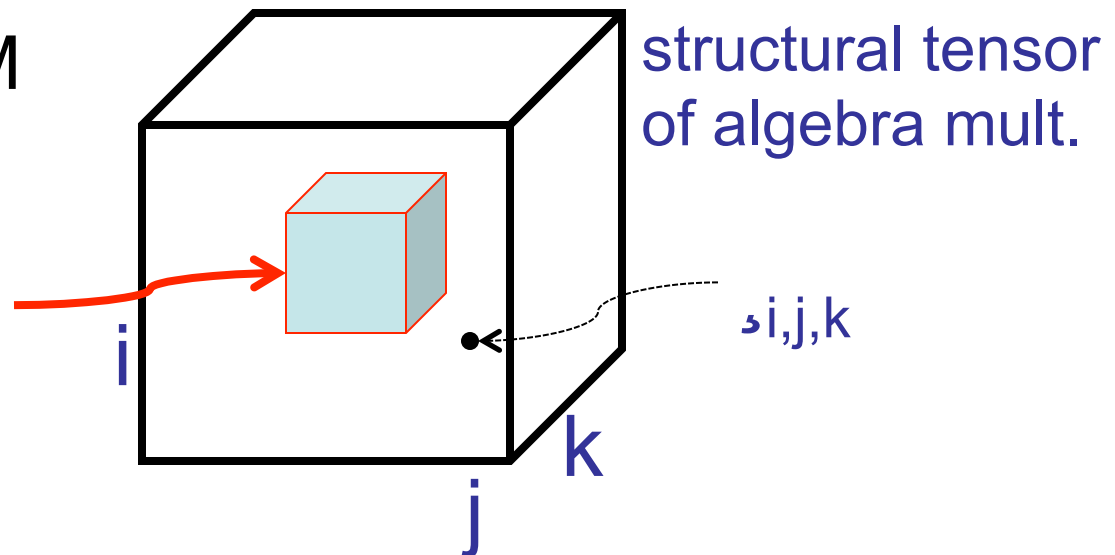  - most ideas foiled by too-large char. degrees

# General semi-simple algebras

- (finite dimensional, complex) algebra specified by
  - "nice basis" $e_1, e_2, \ldots, e_r$
  - structure constants $\flat_{i,j,k}$ satisfying

$$e_i \, e_j = \sum_k \flat_{i,j,k} \, e_k$$

"realizes" MM
if contains*:

MM tensor
$\langle n,n,n \rangle$

structural tensor
of algebra mult.

$\flat_{i,j,k}$

i

j

k

# Weighted vs. unweighted MM

- Technical problem:
  - MM tensor $<n,n,n>$ given by $\sum_{i,j,k} X_{i,j} Y_{j,k} Z_{k,i}$
  - embedding into algebra bounds rank of tensor given by

$$\sum_{i,j,k,\flat i,j,k} X_{i,j} Y_{j,k} Z_{k,I}$$
$$(\text{with } {}_{\flat i,j,k} \neq 0)$$

  - group algebra: ${}_{\flat i,j,k}$ always 0 or 1

# Weighted vs. unweighted MM

s-rank of tensor T: minimum rank of tensor with same **s**upport as T

Does upper bound on s-rank of MM tensor imply upper bound on ordinary rank?

Example:

| $a_{11}$ | $a_{12}$ |
|---|---|
| $a_{21}$ | $a_{22}$ |

X

| $b_{11}$ | $b_{12}$ |
|---|---|
| $b_{21}$ | $b_{22}$ |

=

| $a_{11}b_{11} +$ $a_{12}b_{21}$ | $a_{11}b_{12} +$ $a_{12}b_{22}$ |
|---|---|
| $a_{21}b_{11} +$ $a_{22}b_{21}$ | $a_{21}b_{12} +$ $a_{22}b_{22}$ |

# Weighted vs. unweighted MM

s-rank of tensor T: minimum rank of tensor with same **_s_**upport as T

Does upper bound on s-rank of MM tensor imply upper bound on ordinary rank?

Example:

| $a_{11}$ | $a_{12}$ |
|---|---|
| $a_{21}$ | $a_{22}$ |

X

| $b_{11}$ | $b_{12}$ |
|---|---|
| $b_{21}$ | $b_{22}$ |

!

| $a_{11}b_{11} +$ $a_{12}b_{21}$ | $a_{11}b_{12} +$ $a_{12}b_{22}$ |
|---|---|
| $a_{21}b_{11} +$ $a_{22}b_{21}$ | $a_{21}b_{12} +$ **2¢**$a_{22}b_{22}$ |

does it help if can compute this in
6 multiplications?

# Weighted vs. unweighted MM

- s-rank can be much smaller than rank:

$®$ = n-th root of unity

| 0 | 1 | 1 | 1 |
|---|---|---|---|
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 |
| 1 | 1 | 1 | 0 |

rank n

same support:

| $®_0$ | $®_1$ | $®_2$ | $®_3$ |
|---|---|---|---|
| $®_3$ | $®_0$ | $®_1$ | $®_2$ |
| $®_2$ | $®_3$ | $®_0$ | $®_1$ |
| $®_1$ | $®_2$ | $®_3$ | $®_0$ |

rank 1

-

| 1 | 1 | 1 | 1 |
|---|---|---|---|
| 1 | 1 | 1 | 1 |
| 1 | 1 | 1 | 1 |
| 1 | 1 | 1 | 1 |

rank 1

maybe it's easy to show s-rank of n £ n matrix multiplication is $n^2$ (!!)

# Weighted vs. unweighted MM

$$\omega = \inf \{ \tau : \operatorname{rank}(\langle n,n,n \rangle) \cdot O(n^\tau) \}$$

$$\omega_s = \inf \{ \tau : s\text{-rank}(\langle n,n,n \rangle) \cdot O(n^\tau) \}$$

**Theorem**: $\omega \cdot (3\omega_s - 2)/2$

in particular, $\omega_s \cdot 2 + \varepsilon \Rightarrow \omega \cdot 2 + (3/2)\varepsilon$

- Proof idea:
  – find ¼ $n^2$ copies of $\langle n,n,n \rangle$ in 3$^{\text{rd}}$ tensor power
  – when broken up this way, can rescale

# A promising family of semisimple algebras

# Coherent configurations

"group theory without groups"

- points $X$, partition $R_1, R_2, ..., R_r$ of $X^2$
  - diagonal $\{(x,x) : x \in X\}$ is union of some classes
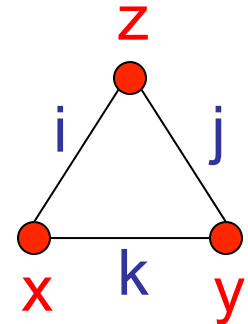
  if *one* class: "association scheme"

  - for each i, there is i*

  $p_{i,j}{}^k = p_{j,i}{}^k$ : commutative

  $R_{i*} = \{(y,x) : (x,y) \in R_i\}$

  - exist integers $p_{i,j}{}^k$ such that for all $(x,y) \in R_k$:

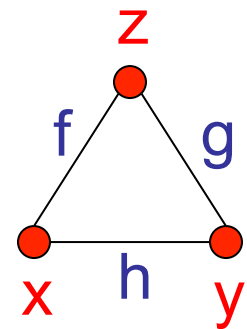  $\#\{z: (x,z) \in R_i$ and $(z,y) \in R_j\} = p_{i,j}{}^k$

# Coherent configs: examples

- Hamming scheme:
  - points 0/1 vectors
  - classes determined by hamming distance


- distance-regular graph:
  - points = vertices
  - classes determined by distance in graph metric

# Coherent configs: examples

- scheme based on finite group G
  - set X = finite group G
  - classes $R_g = \{(x, xg) : x \in X\}$

$$p_{f,g}^{\ h} = 1 \text{ if } fg = h, 0 \text{ otherwise}$$

- "Schurian":
  - group G acts on set X
  - classes = orbits of (diagonal) G-action on $X^2$

# Coherent configs: examples

- "Schurian":
  - group G acts on set X
  - classes = orbits of (diagonal) G-action on $X^2$

- one Schurian scheme: "group scheme"
  - group G x G acts on G via $(g,h)¢x = gxh^{-1}$
  - orbits all of the form {(x,y): $xy^{-1}$ 2 $C_i$} for conjugacy class $C_i$
  - always commutative!

# Adjacency algebra

CC: points $X$, partition $R_1, R_2, \ldots, R_r$ of $X^2$

- for each class $R_i$, matrix $A_i$ with

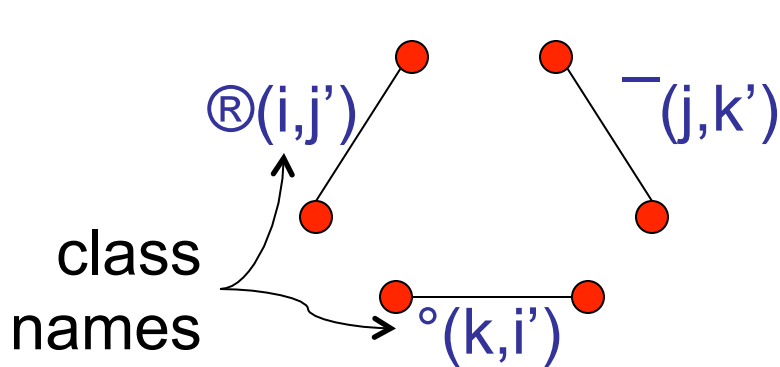$$A_i[x,y] = 1 \text{ iff } (x,y) \, 2 \, R_i$$

- 3 CC axioms )
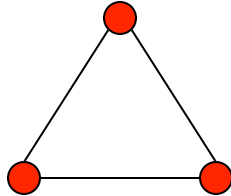
$\{A_i\}$ generate a semisimple algebra

– e.g., $3^{rd}$ axiom implies $A_i A_j = \sum_k p_{ij}^k A_k$

– if the CC based on group $G$, algebra is $C[G]$

# Nice basis conditions

- group algebra C[G]: "nice basis" yields triple product property

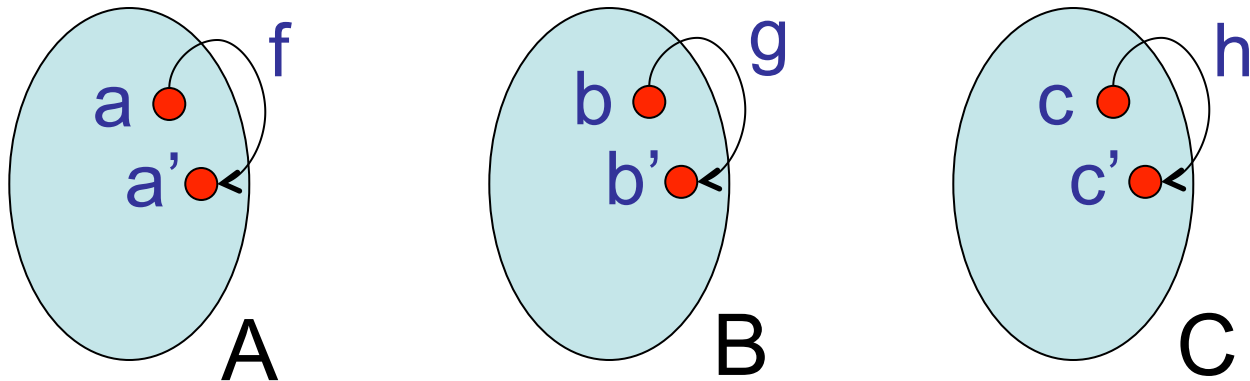- adjacency algebras of CCs: "nice basis" yields triangle condition:

®(i,j')     ¯(j,k')     can look like

class names     °(k,i')     iff i = i', j = j', k = k'

# Nice basis conditions

- Schurian CCs: "nice basis" yields
  - group G acts on set X
  - subsets A,B,C of X realize <|A|, |B|, |C|> if:



fgh = 1 implies a = a', b = b', c = c'

# Coherent configs vs. groups

Generalization for generalization's sake?

- recall group framework:

  – non-commutative necessary

**Theorem**: in group G realizing n£n matrix multiplication, with character degrees $d_1$, $d_2$, $d_3$,…, we obtain:

$$R(<n,n,n>) \cdot \sum_i d_i^{\omega} \cdot d_{max}^{\omega-2} \not| |G|$$

goals: $|G|$ ¼ $n^2$ *and* small $d_{max}$

# Coherent configs vs. groups

Generalization for generalization's sake?

- coherent configuration framework:
  - commutative suffices!

  - combinatorial constructions from old setting yield
    $$!_s < 2.48, \; !_s < 2.41$$
  - conjectures from old setting (if true) would imply $!_s = 2$

in commutative Schurian CC's

even group schemes

even symmetric

Sept. 17, 2014

# Proof idea

we prove a general transformation:

if can realize <span style="color:red">several independent matrix multiplications</span> in CC…

- can do this in abelian groups
- conjectures: can "pack optimally"

… then high <span style="color:red">symmetric power</span> of CC realizes *single* matrix multiplication

– reproves Schönhage's

Asymptotic Sum Inequality

preserves commutativity

# Commutative CCs suffice

**Main point**

embedding n x n matrix multiplication into a commutative coherent config-uration of rank ¼ $n^2$ is a viable route to ! = 2

(no representation theory needed)

# Open problems

- find a construction in new framework that
  - proves non-trivial bound on $!_s$
  - is not based on constructions from old setting

- is the (border) s-rank of <2,2,2> = 6?

- embed n £ n MM into commutative coherent configuration of rank ¼ $n^2$