



Interactive Combinatorial Supercomputing

John R. Gilbert

University of California, Santa Barbara

Viral Shah, Imran Patel (UCSB)

Alan Edelman (MIT and Interactive Supercomputing)

Ron Choy, David Cheng (MIT)

Parry Husbands (Lawrence Berkeley Lab)

Steve Reinhardt, Todd Letsche (SGI)

Support: DOE Office of Science, DARPA, SGI, ISC

Parallel Computing Today



Columbia,
NASA Ames Research Center

Departmental Beowulf cluster



But!

How do you program it?

C with MPI

```
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#include "mpi.h"
#define A(i,j) ( 1.0/(((1.0*(i)+(j))*(1.0*(i)+(j)+1)/2 + (1.0*(i)+1)) )
void errorExit(void);
double normalize(double* x, int mat_size);
int main(int argc, char **argv)
{
    int num_procs;
    int rank;
    int mat_size = 64000;
    int num_components;
    double *x = NULL;
    double *y_local = NULL;
    double norm_old = 1;
    double norm = 0;
    int i,j;
    int count;
    if (MPI_SUCCESS != MPI_Init(&argc, &argv)) exit(1);
    if (MPI_SUCCESS != MPI_Comm_size(MPI_COMM_WORLD,&num_procs)) errorExit();
```

C with MPI (2)

```
if (0 == mat_size % num_procs) num_components = mat_size/num_procs;
else num_components = (mat_size/num_procs + 1);
mat_size = num_components * num_procs;
if (0 == rank) printf("Matrix Size = %d\n", mat_size);
if (0 == rank) printf("Num Components = %d\n", num_components);
if (0 == rank) printf("Num Processes = %d\n", num_procs);
x = (double*) malloc(mat_size * sizeof(double));
y_local = (double*) malloc(num_components * sizeof(double));
if ( (NULL == x) || (NULL == y_local) )
{
    free(x);
    free(y_local);
    errorExit();
}
if (0 == rank)
{
    for (i=0; i<mat_size; i++)
    {
        x[i] = rand();
    }
    norm = normalize(x,mat_size);
}
```

C with MPI (3)

```
if (MPI_SUCCESS !=
    MPI_Bcast(x, mat_size, MPI_DOUBLE, 0, MPI_COMM_WORLD)) errorExit();
count = 0;
while (fabs(norm-norm_old) > TOL) {
    count++;
    norm_old = norm;
    for (i=0; i<num_components; i++)
    {
        y_local[i] = 0;
    }
    for (i=0; i<num_components && (i+num_components*rank)<mat_size; i++)
    {
        for (j=mat_size-1; j>=0; j--)
        {
            y_local[i] += A(i+rank*num_components,j) * x[j];
        }
    }
    if (MPI_SUCCESS != MPI_Allgather(y_local, num_components, MPI_DOUBLE, x,
        num_components, MPI_DOUBLE, MPI_COMM_WORLD)) errorExit();
}
```

C with MPI (4)

```
    norm = normalize(x, mat_size);
}
if (0 == rank)
{
    printf("result = %16.15e\n", norm);
}
free(x);
free(y_local);

MPI_Finalize();
exit(0);
}
void errorExit(void)
{
    int rank;
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    printf("%d died\n", rank);
    MPI_Finalize();
    exit(1);
}
```

C with MPI (5)

```
double normalize(double* x, int mat_size)
{
    int i;
    double norm = 0;
    for (i=mat_size-1; i>=0; i--)
    {
        norm += x[i] * x[i];
    }
    norm = sqrt(norm);
    for (i=0; i<mat_size; i++)
    {
        x[i] /= norm;
    }
    return norm;
}
```



```
A = rand(4000*p, 4000*p);  
x = randn(4000*p, 1);  
y = zeros(size(x));  
while norm(x-y) / norm(x) > 1e-11  
    y = x;  
    x = A*x;  
    x = x / norm(x);  
end;
```

Background

- Matlab*P 1.0 (1998): Edelman, Husbands, Isbell (MIT)
- Matlab*P 2.0 (2002-): MIT / UCSB / LBNL
- Star-P (2004-): Interactive Supercomputing / SGI



Data-Parallel Operations

< M A T L A B >

Copyright 1984-2001 The MathWorks, Inc.

Version 6.1.0.1989a Release 12.1

```
>> A = randn(500*p, 500*p)
```

```
A = ddense object: 500-by-500
```

```
>> E = eig(A);
```

```
>> E(1)
```

```
ans = -4.6711 +22.1882i
```

```
e = pp2matlab(E);
```

```
>> ppwhos
```

Name	Size	Bytes	Class
A	500px500p	688	ddense object
E	500px1	652	ddense object
e	500x1	8000	double array (complex)

Task-Parallel Operations

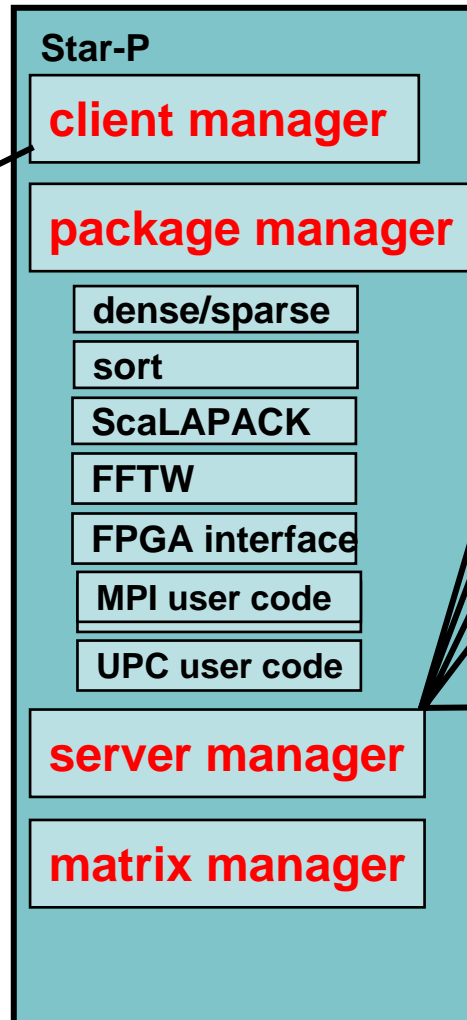
```
>> quad('4./(1+x.^2)', 0, 1);  
ans = 3.14159270703219  
  
>> a = (0:3*p) / 4  
a = ddense object: 1-by-4  
  
>> a(:, :)  
ans =  
  
0  
0.2500000000000000  
0.5000000000000000  
0.7500000000000000  
  
>> b = a + .25;  
  
>> c = ppeval('quad', '4./(1+x.^2)', a, b);  
c = ddense object: 1-by-4  
  
>> sum(c)  
ans = 3.14159265358979
```

Star-P Architecture

MATLAB®



Ordinary Matlab variables



processor #0

processor #1

processor #2

processor #3

⋮

processor #n-1

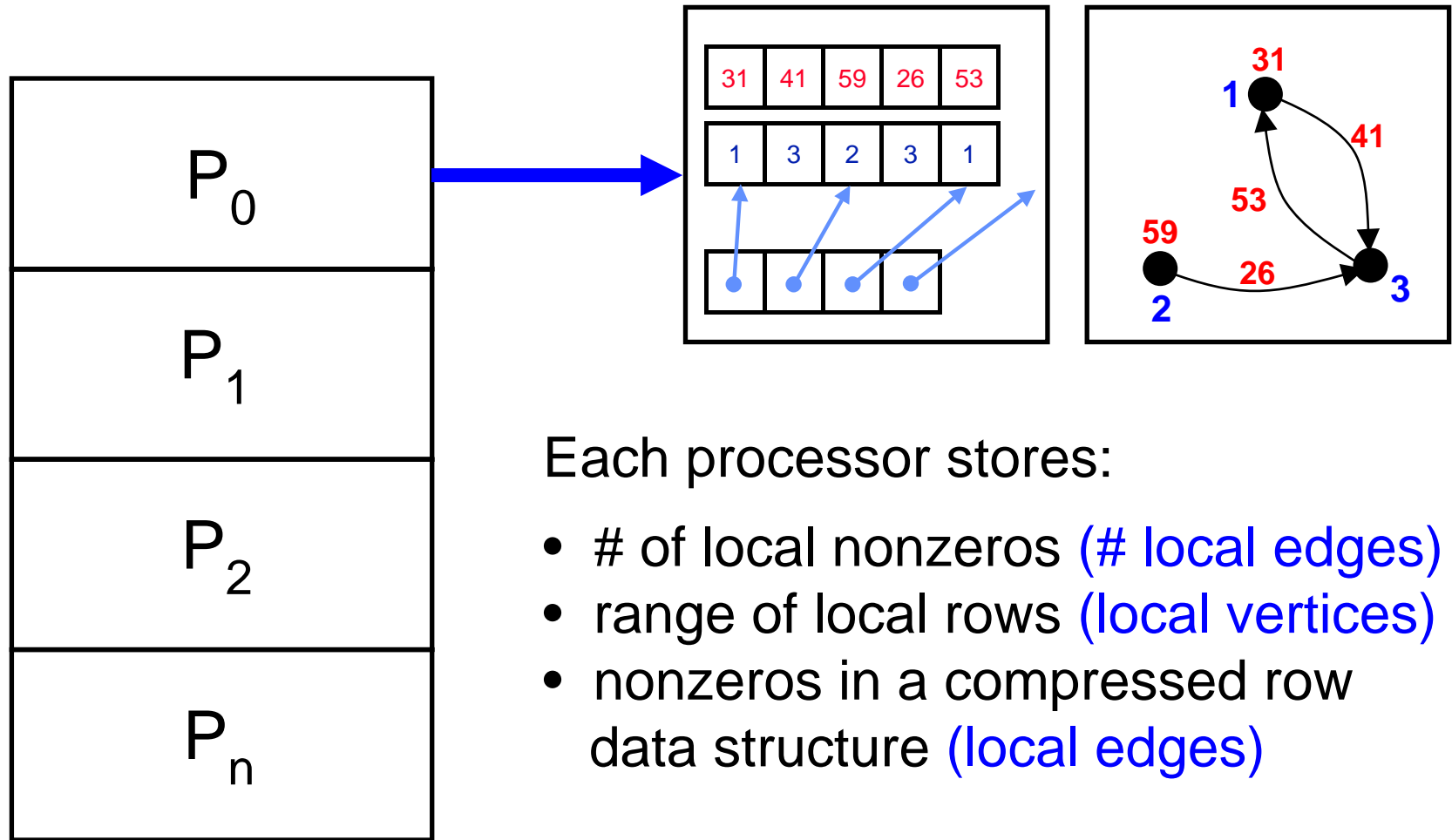
Distributed matrices

Matlab sparse matrix design principles

- All operations should give the same results for sparse and full matrices (almost all)
- Sparse matrices are never created automatically, but once created they propagate
- Performance is important -- but usability, simplicity, completeness, and robustness are more important
- Storage for a sparse matrix should be $O(\text{nonzeros})$
- Time for a sparse operation should be $O(\text{flops})$ (as nearly as possible)

**Star-P dsparse matrices: same principles,
but some different tradeoffs**

Distributed sparse array structure



The sparse() constructor

- $A = \text{sparse}(I, J, V, nr, nc);$
- Input: dense vectors I, J, V , dimensions nr, nc
- Output: $A(I(k), J(k)) = V(k)$
- Sum values with duplicate indices
- Sorts triples $\langle i, j, v \rangle$ by $\langle i, j \rangle$
- Inverse: $[I, J, V] = \text{find}(A);$

Sparse array and matrix operations

- `dsparse` layout, same semantics as `ddense`
- Matrix arithmetic: `+`, `max`, `sum`, etc.
- `matrix * matrix` and `matrix * vector`
- Matrix indexing and concatenation
$$A(1:3, [4\ 5\ 2]) = [B(:, J)\ C];$$
- Linear solvers: `x = A \ b`; using SuperLU (MPI)
- Eigensolvers: `[V, D] = eigs(A)`; using PARPACK (MPI)

Sparse matrix times dense vector

- $y = A * x$
- First matvec with A caches a communication schedule
- Later matvecs with A use the cached schedule
- Communication and computation overlap

Combinatorial Scientific Computing

- Sparse matrix methods
- Knowledge discovery
- Web search and information retrieval
- Graph matching
- Machine learning
- Geometric modeling
- Computational biology
- Bioinformatics
- . . .

How will combinatorial methods be used by nonexperts?

Analogy: Matrix division in Matlab

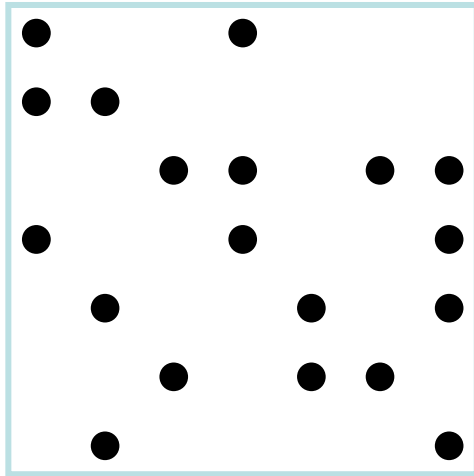
$$\mathbf{x} = \mathbf{A} \setminus \mathbf{b};$$

- Works for either full or sparse A
- Is A square?
 - no** => use QR to solve least squares problem
- Is A triangular or permuted triangular?
 - yes** => sparse triangular solve
- Is A symmetric with positive diagonal elements?
 - yes** => attempt Cholesky after symmetric minimum degree
- Otherwise
 - =>** use LU on $\mathbf{A}(:, \text{colamd}(\mathbf{A}))$

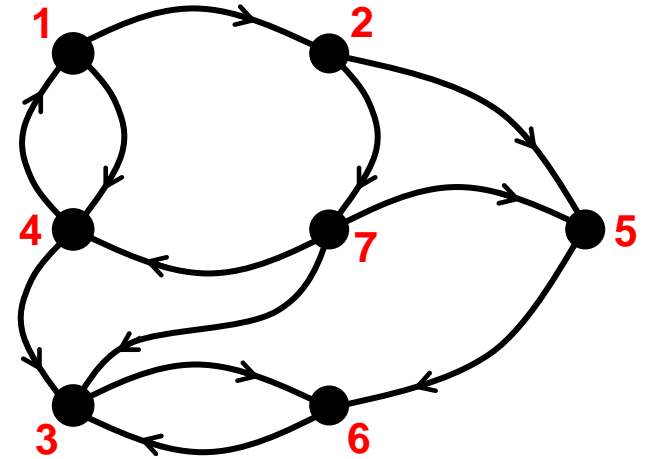
Combinatorics in Star-P

- Represent a graph as a sparse adjacency matrix
- A sparse matrix language is a good start on primitives for computing with graphs
 - Random-access indexing: `A(i, j)`
 - Neighbor sequencing: `find (A(i, :))`
 - Sparse table construction: `sparse (I, J, V)`
 - Breadth-first search step : `A * v`

Sparse adjacency matrix and graph

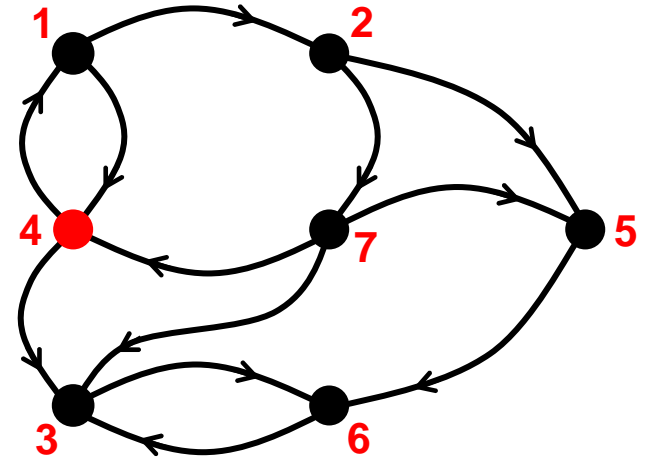
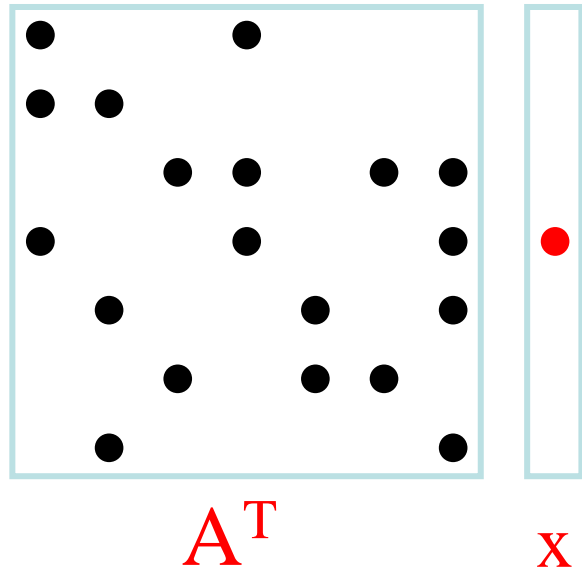


A^T



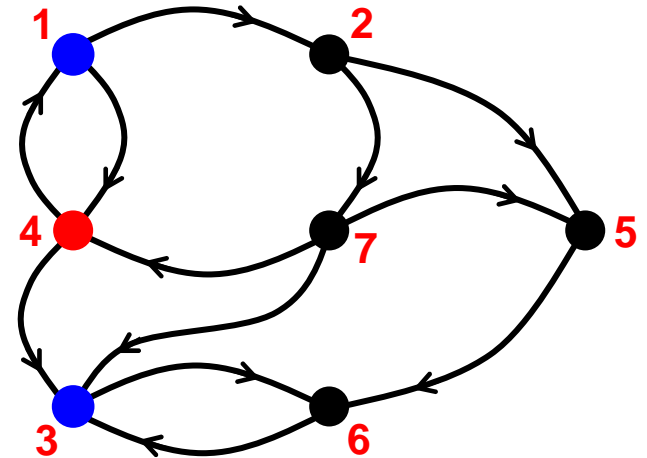
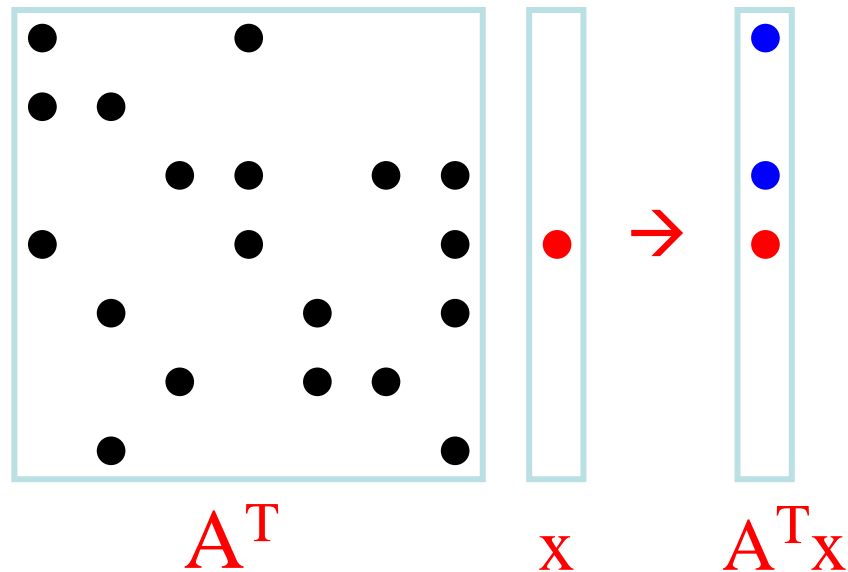
- Adjacency matrix: sparse array w/ nonzeros for graph edges
- Storage-efficient implementation from sparse data structures

Breadth-first search: sparse mat * vec



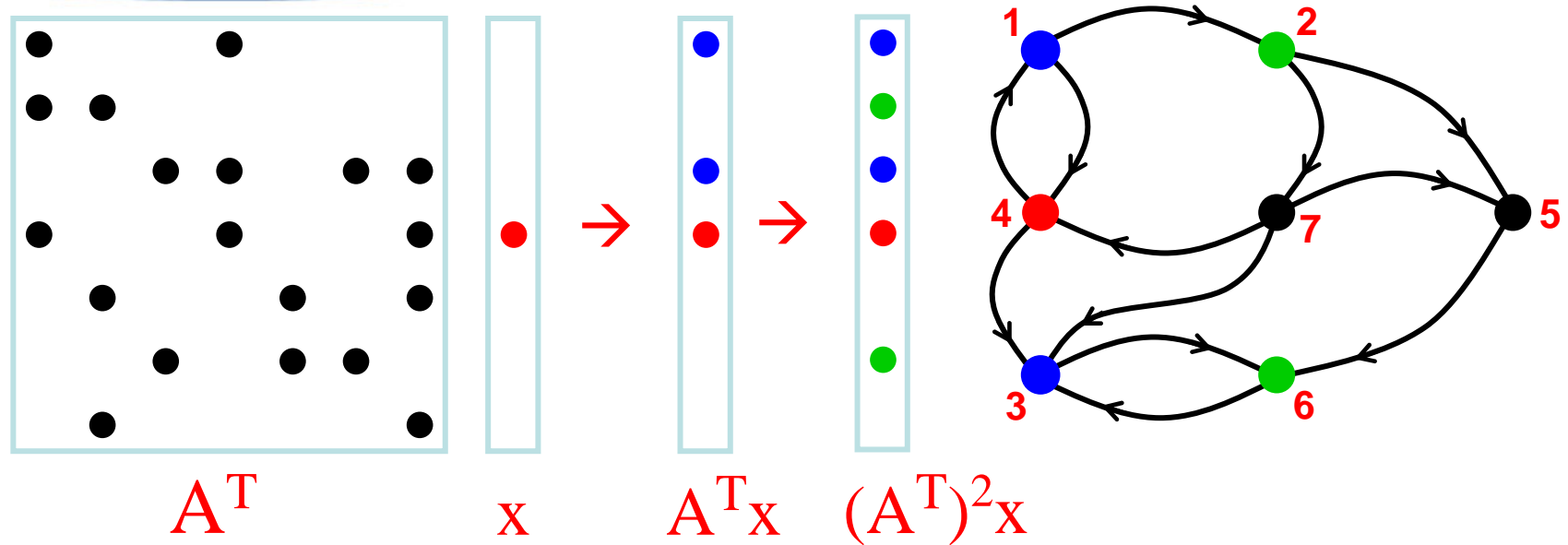
- Multiply by adjacency matrix \rightarrow step to neighbor vertices
- Efficient implementation from sparse data structures

Breadth-first search: sparse mat * vec



- Multiply by adjacency matrix \rightarrow step to neighbor vertices
- Efficient implementation from sparse data structures

Breadth-first search: sparse mat * vec



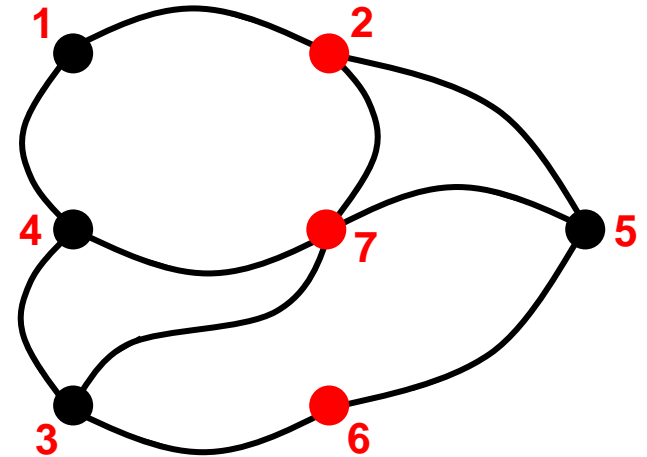
- Multiply by adjacency matrix \rightarrow step to neighbor vertices
- Efficient implementation from sparse data structures

Connected components of a graph

- Sequential Matlab uses depth-first search (**dmperm**), which doesn't parallelize well
- Pointer-jumping algorithms (Shiloach/Vishkin & descendants)
 - repeat
 - Link every (super)vertex to a neighbor
 - Shrink each tree to a supervertex by pointer jumping
 - until no further change
- Other coming graph kernels:
 - Shortest-path search (after Husbands, LBNL)
 - Bipartite matching (after Riedy, UCB)
 - Strongly connected components (after Pinar, LBNL)

Maximal independent set

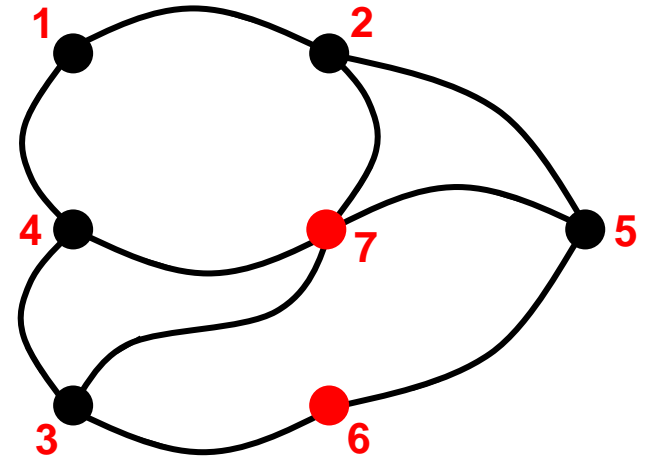
```
degree = sum(G, 2);  
prob = 1 ./ (2 * deg);  
select = rand (n, 1) < prob;  
  
if ~isempty (select & (G * select));  
    % keep higher degree vertices  
end  
  
IndepSet = [IndepSet select];  
  
neighbor = neighbor | (G * select);  
remain = neighbor == 0;  
G = G(remain, remain);
```



Starting guess:
Select some vertices
randomly

Maximal independent set

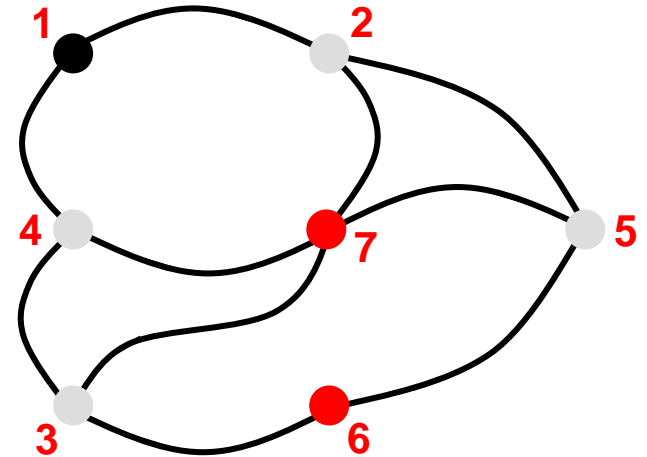
```
degree = sum(G, 2);  
prob = 1 ./ (2 * deg);  
select = rand (n, 1) < prob;  
  
if ~isempty (select & (G * select))  
    % keep higher degree vertices  
end  
  
IndepSet = [IndepSet select];  
  
neighbor = neighbor | (G * select);  
remain = neighbor == 0;  
G = G(remain, remain);
```



If neighbors are selected, keep only a higher-degree one.
Add selected vertices to the independent set.

Maximal independent set

```
degree = sum(G, 2);  
prob = 1 ./ (2 * deg);  
select = rand (n, 1) < prob;  
  
if ~isempty (select & (G * select));  
    % keep higher degree vertices  
end  
  
IndepSet = [IndepSet select];  
  
neighbor = neighbor | (G * select);  
remain = neighbor == 0;  
G = G(remain, remain);
```

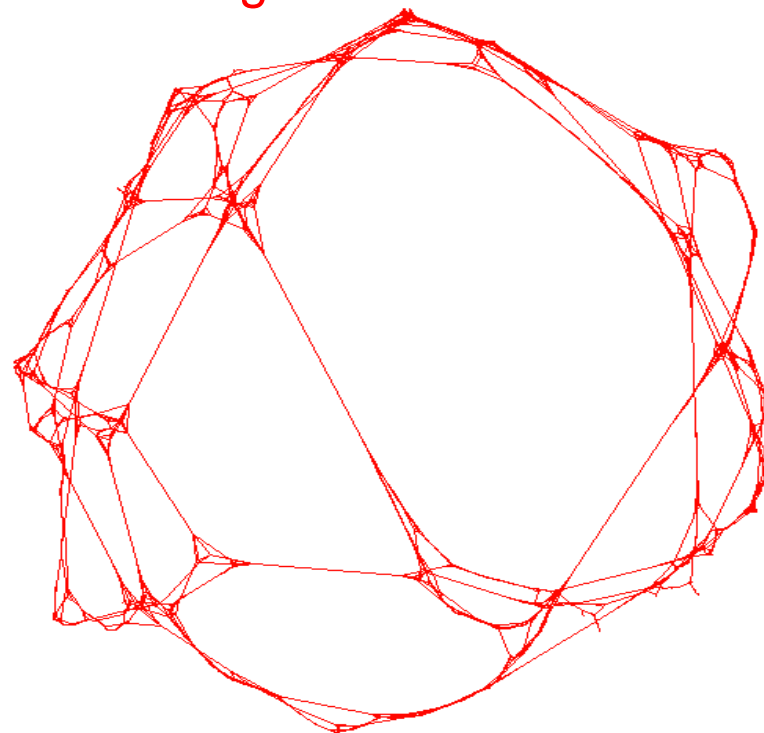
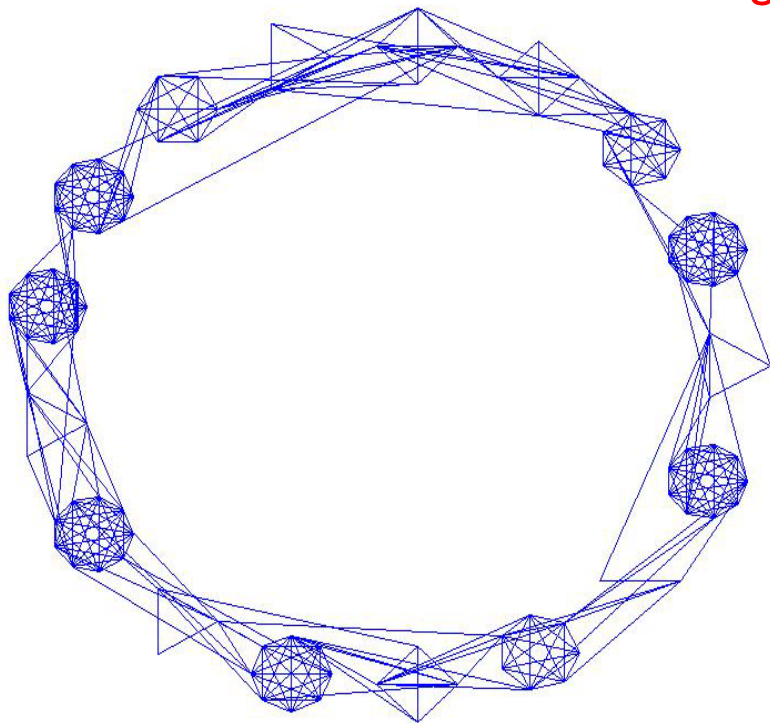


Discard neighbors of
the independent set.

Iterate on the rest of
the graph.

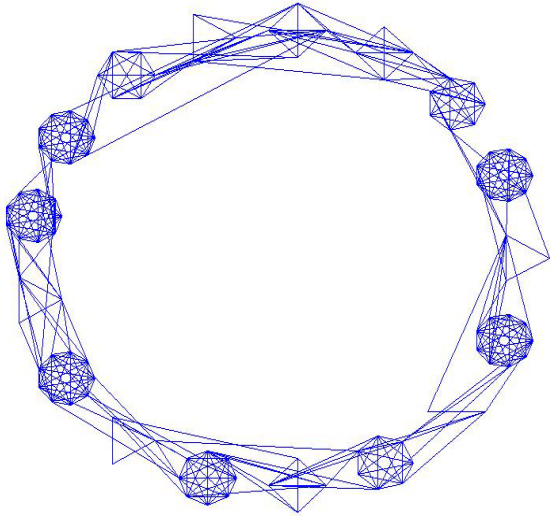
SSCA#2: “Graph Analysis”

Fine-grained, irregular data access
Searching and clustering



- Many tight clusters, loosely interconnected
- Input data is edge triples $\langle i, j, \text{label}(i,j) \rangle$
- Vertices and edges permuted randomly

SSCA#2: Graph statistics

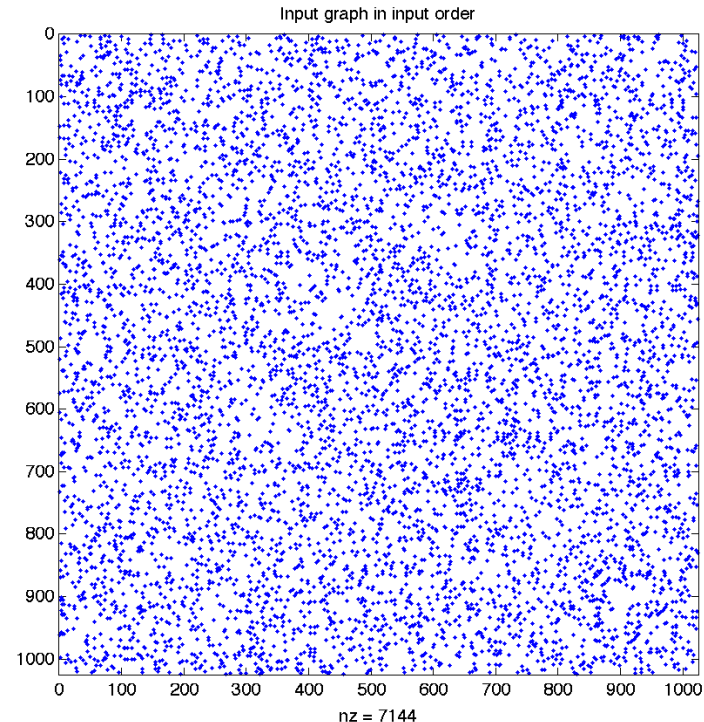


- Scalable data generator
- Given “scale” = $\log_2(\text{\#vertices})$
- Creates edge triples $\langle i, j, \text{label}(i,j) \rangle$
- Randomly permutes triples and vertex numbers

Scale	#Vertices	#Cliques	#Edges Directed	#Edges Undirected
10	1,024	186	13,212	3,670
15	32,768	2,020	1,238,815	344,116
20	1,048,576	20,643	126,188,649	35,052,403
25	33,554,432	207,082	12,951,350,000	3,597,598,000
30	1,073,741,824	2,096,264	1,317,613,000,000	366,003,600,000

Statistics for SSCA2 spec v1.1

Concise SSCA#2 in Star-P



Kernel 1: Construct graph data structures

- Graphs are dsparse matrices, created by `sparse()`

Kernels 2 and 3

Kernel 2: Search by edge labels

- About 12 lines of executable Matlab or Star-P

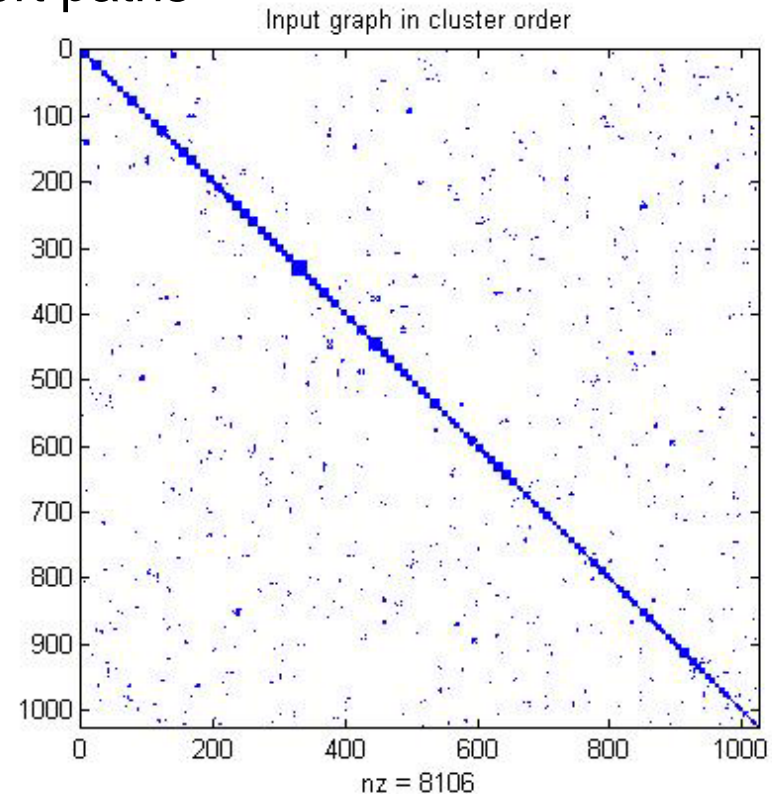
Kernel 3: Extract subgraphs

- Returns subgraphs consisting of vertices and edges within fixed distance of given starting vertices
- Sparse matrix-matrix product for multiple breadth-first search
- About 25 lines of executable Matlab or Star-P

Kernel 4: Clustering by BFS

- Grow local clusters from many seeds in parallel
- Breadth-first search by sparse matrix * matrix
- Cluster vertices connected by many short paths

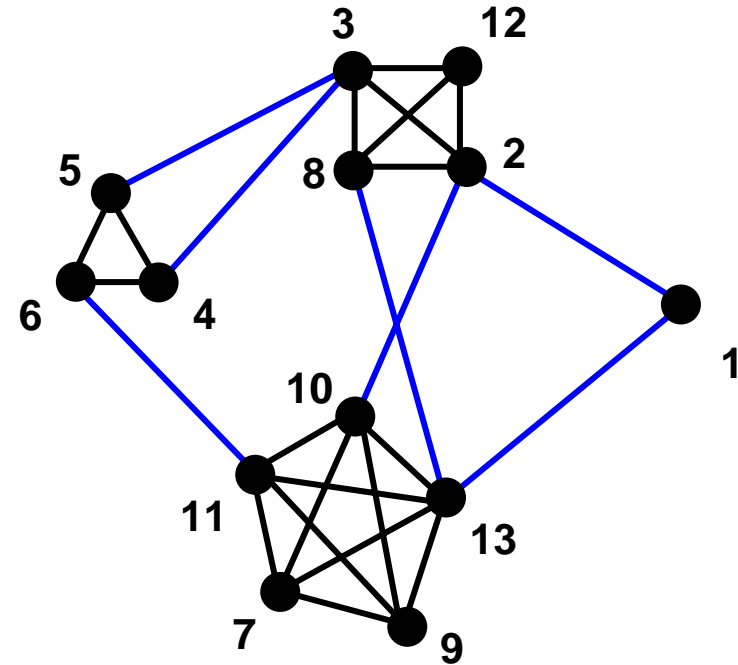
```
% Grow each seed to vertices  
%   reached by at least k  
%   paths of length 1 or 2  
  
C = sparse(seeds, 1:ns, 1, n, ns);  
C = A * C;  
C = C + A * C;  
C = C >= k;
```



Kernel 4: Clustering by peer pressure

Steps in a peer pressure algorithm:

1. Vote for a cluster leader
2. Collect neighbor votes
3. Vote for a new leader
(based on neighbor votes)



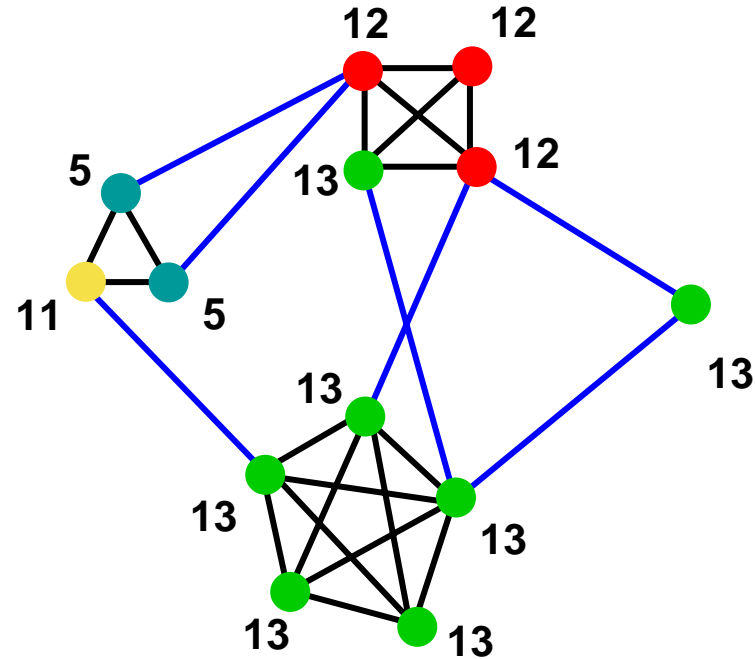
- Clustering qualities depend on details of each step.
- Want relatively few potential leaders, e.g. a maximal indep set. Other choices possible – for SSCA2 graph, simpler rules work too.
- Neighbor votes can be combined using various weightings.
- Each version of kernel4 is about 25 lines of code.

Kernel 4: Clustering by peer pressure

```
[ignore, leader] = max(G);
```

```
S = G * sparse(1:n, leader, 1, n, n);
```

```
[ignore, leader] = max(S);
```



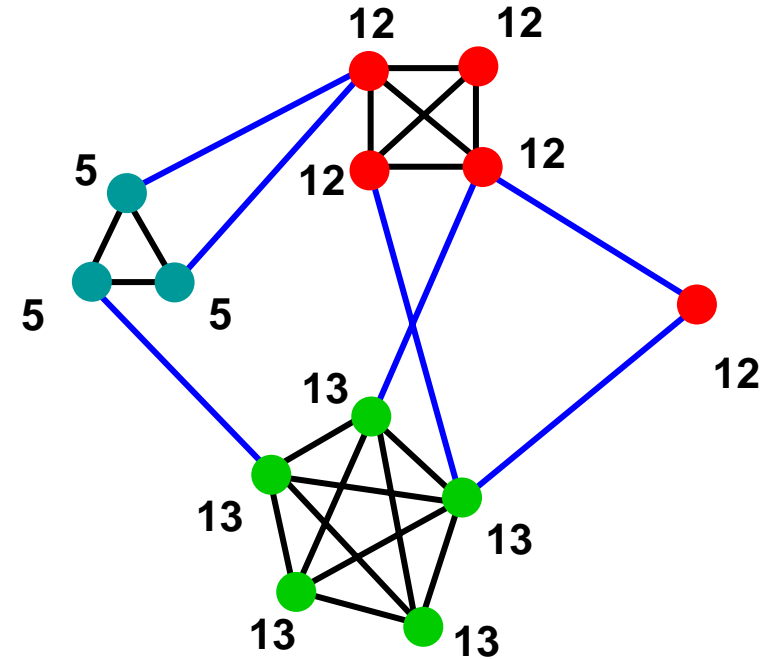
- Each vertex votes for highest numbered neighbor as its leader
- Number of leaders is approximately number of clusters (small relative to the number of nodes)

Kernel 4: Clustering by peer pressure

```
[ignore, leader] = max(G);
```

```
S = sparse(leader, 1:n, 1, n, n) * G;
```

```
[ignore, leader] = max(S);
```



- Matrix multiplication gathers neighbor votes
- $S(i,j)$ is # of votes for i from j 's neighbors
- In SSCA2 (spec1.0), most of graph structure is recovered right away; iteration needed for harder graphs

Expressive Power: SSCA#2 Kernel 3

Star-P (25 lines)

```
A = spones(G.edgeWeights{1});
nv = max(size(A));
npar = length(G.edgeWeights);
nstarts = length(starts);
for i = 1:nstarts
    v = starts(i);
    % x will be a vector whose nonzeros
    % are the vertices reached so far
    x = zeros(nv,1);
    x(v) = 1;
    for k = 1:pathlen
        x = A*x;
        x = (x ~= 0);
    end;
    vtxmap = find(x);
    S.edgeWeights{1} = G.edgeWeights{1}(vtxmap,vtxmap);
    for j = 2:npar
        sg = G.edgeWeights{j}(vtxmap,vtxmap);
        if nnz(sg) == 0
            break;
        end;
        S.edgeWeights{j} = sg;
    end;
    S.vtxmap = vtxmap;
    subgraphs{i} = S;
end
```

MATLABmpi (91 lines)

```
declareGlobals;

% Wait for a response for each request we sent out.
for unused = 1:numpar
    [src tag] = probeSubgraphs(G, [P tag K3 dataReq]);
    [starts newEdges] = MPI_Recv(src, tag, P.comm);
    subg.edgeWeights{1}(:, starts) = newEdges;
    [newEnds unused] = find(newEdges);
    allNewEnds = [allNewEnds; newEnds];
end
end % of if -P paral

% Eliminate any new ends already in the all starts list.
newStarts = setdiff(allNewEnds, allStarts);
allStarts = [allStarts; newStarts];

% ENABLE_PLOT_KIDB
plotEdges(subg.edgeWeights{1}, startVertex, endVertex, k);
end % of ENABLE_PLOT_KIDB

if isempty(newStarts) % if empty we can quit early.
    break;
end
end
```

<u>Lines of code</u>	Star-P cSSCA2	MATLABmpi spec	C/Pthreads/ SIMPLE
Kernel 1	29	68	256
Kernel 2	12	44	121
Kernel 3	25	91	297
Kernel 4	44	295	241

```
% For each processor which has any of the vertices we need.
startDests = floor((newStarts - 1) / P.myV);
uniqDests = unique(startDests);
for dest = uniqDests
    starts = newStarts(startDests == dest);

    if dest == P.myRank
        newEdges = G.edgeWeights{1}(:, starts - P.myBase);
        subg.edgeWeights{1}(:, starts) = newEdges;
        [allNewEnds unused] = find(newEdges);
    else
        % empty starts
        MPI_Send(dest, P.tag K3 dataReq, P.comm, starts);
        numRecv = numRecv + 1;
    end
end
```

Scaling up

Recent results on SGI Altix (up to 128 processors):

- Have run SSCA2 on graphs with $2^{27} = 134$ million vertices and about one billion (10^9) edges (spec v1.0)
- Benchmarking in progress for spec v1.1 (different graph generator)
- Have manipulated graphs with 400 million vertices and 4 billion edges
- Timings scale well – for large graphs,
 - 2x problem size \rightarrow 2x time
 - 2x problem size & 2x processors \rightarrow same time

Using this benchmark to tune lots of infrastructure

Work in progress: Toolbox for Graph Analysis and Pattern Discovery

Layer 1: Graph Theoretic Tools

- Graph operations
- Global structure of graphs
- Graph partitioning and clustering
- Graph generators
- Visualization and graphics
- Scan and combining operations
- Utilities

