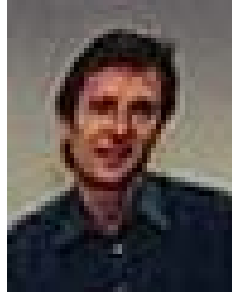


Efficient Statistical Learning from “Big Science” Data



Andrew
Moore
Auton Lab



Andrew
Connolly
U. Pittsburgh



Jeremy
Kubica
Auton Lab



Ting
Liu
Auton Lab



The Auton Lab
School of Computer Science
Carnegie Mellon University
www.autonlab.org

The Auton Lab

Faculty:	Andrew Moore (Prof), Jeff Schneider (Research Scientist), Artur Dubrawski (Systems Scientist)
Postdoctoral Fellows:	Brigham Anderson, Alexander Gray, Paul Komarek, Dan Pelleg
Graduate Students:	Brent Bryan, Kaustav Das, Khalid El-Arini, Anna Goldenberg, Jeremy Kubica, Ting Liu, Daniel Neill, Sajid Siddiqi, Purna Sarkar, Ajit Singh, Weng-Keen Wong
Head of Software Development:	Jeanie Komarek
Programmers:	Patrick Choi, Adam Goode, Pat Gunn, Joey Liang, John Ostlund, Robin Sabhnani, Rahul Sankathar
Executive Assistant:	Kristen Schrauder
Head of Sys. Admin:	Jacob Joseph
Undergraduate and Masters Interns:	Kenny Daniel, Sandy Hsu, Dongryeol Lee, Jennifer Lee, Avilay Parekh, Chris Rotella, Jonathan Terleski
Recent Alumni:	Drew Bagnell (RI faculty), Scott Davies (Google), David Cohn (Google), Geoff Gordon (CMU), Paul Hsiung (USC), Marina Meila (U. Washington), Remi Munos (Ecole Polytechnique), Malcolm Strens (Qinetiq)

Current Sponsors

- National Science Foundation (NSF)
- NASA
- Defense Advanced Research Projects Agency (DARPA)
- Department of Homeland Security (DHS)
- Homeland Security Advanced Research Projects Agency (HSARPA)
- United States Department of Agriculture (USDA)
- State of Pennsylvania
- Pfizer Corporation
- Caterpillar Corporation
- British Petroleum
- Psychogenics Corporation
- Transform Pharmaceuticals
- Health Canada

Collaborators...



Alex Gray



Daniel Neill



Paul Komarek



Ting Liu



Kan Deng

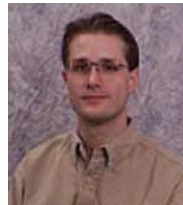


Jeremy Kubica

Drew Bagnell



Paul Hsiung



Brigham Anderson



Scott Davies



Weng-Keen Wong



Jeff Schneider



Anna Goldenberg

Ajit Singh



Dan Pelleg



Bob Nichol

Chris Genovese

Larry Wasserman

The "Statistical Data Mining for Data Mining" team: (clockwise from far left) Schneider, Andy Connolly, Larry Wasserman and Chris

Istvan Szapudi

Alex Szalay

Andy Connolly

Greg Cooper

Mike Wagner

Andrew Lawson

Rich Tsui

Sajid Siddiqi

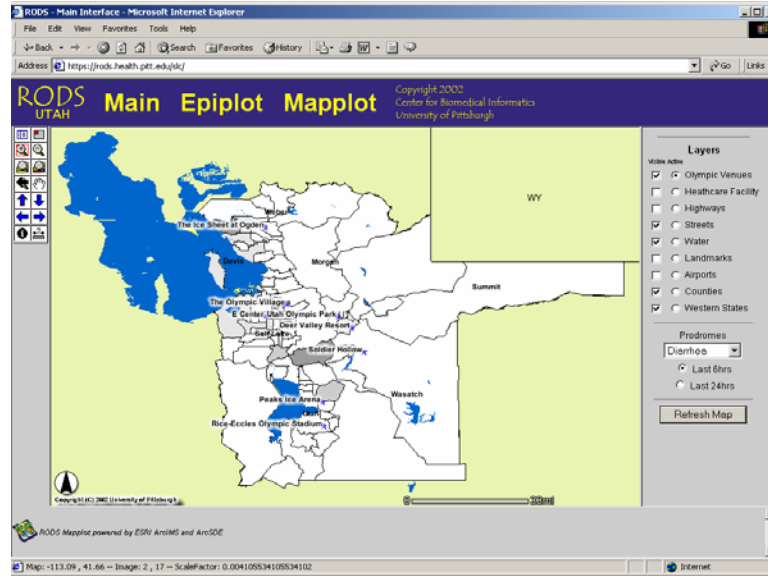
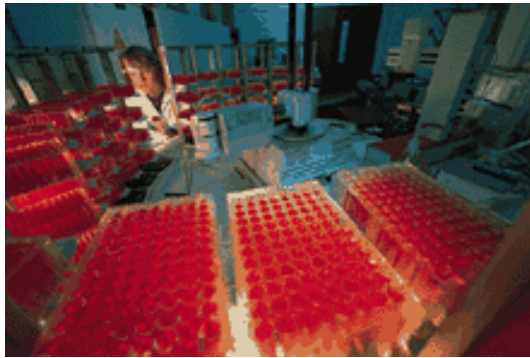
1996	1997	1998	1999	2000	2001	2002	2003
<p>M&M Mars</p> <p>Line control</p> <p>Adrenaline (NOX minimization)</p>	<p>Kodak (Image stabilization)</p> <p>Digital Equipment (pregnancy monitoring)</p>	<p>M&M Mars (manufacturing)</p> <p>NASA/NSF (Astrophysics mining)</p> <p>3M Textile tension control</p>	<p>Caterpillar (Spare parts)</p> <p>US Army (biotoxin detection)</p> <p>M&M Mars: Scheduling with uncertainty</p> <p>3M (Adhesive design)</p>	<p>DigitalMC (Music tastes)</p> <p>Caterpillar (emissions)</p> <p>SmartMoney (anomalies)</p> <p>Unilever (Brand Management)</p> <p>Phillips Petroleum (work-force optimization)</p> <p>Cellomics (screened anomaly detection)</p>	<p>Biometrics company (health monitor)</p> <p>Boeing (intrusion)</p> <p>Masterfoods (new product development)</p> <p>Cellomics (pro-teomics screen)</p> <p>ABB (Circuit-breaker supply chain)</p> <p>SwissAir (Flight delays)</p> <p>3M (secret)</p> <p>Washington Public Hospital System (ER delays)</p> <p>Unilever (targeted marketing)</p>	<p>NASA (National Virtual Observatory)</p> <p>NSF (astrostatistics software)</p> <p>DARPA (national disease monitor)</p> <p>Masterfoods (bio-chemistry)</p> <p>Pfizer (High-throughput screen)</p> <p>Caterpillar Inc. (Self-optimizing Engines)</p> <p>Beverage Company (Ingredients/Manufacturing/Marketing/Sales Bayes Net)</p> <p>Transform Pharma (massive autonomous experiment design)</p> <p>Census Bureau (privacy protection)</p> <p>Psychogenics Inc: Effects of psychotropic drugs on rats</p>	<p>NSF (astrostatistics software)</p> <p>Masterfoods (bio-chemistry)</p> <p>State of PA (National Disease Monitor [with Mike Wagner of U. Pitt])</p> <p>State of PA (Anti Cancer [collaboration with CMU Biology])</p> <p>DARPA (detecting patterns in links)</p> <p>Other Government Departments (identifying dangerous people, potential collaborators, and aliases)</p> <p>Other Government Departments (detecting a class of clusters)</p> <p>Other Pharma Research Co. Life Science specific data mining</p> <p>United States Department of Agriculture: Early warning system for food terrorism</p> <p>NSF: Biosurveillance Algorithms</p>

Auton/SPR Deployments

Our 5 biggest applications in 2004

Biomedical Security (with Mike Wagner, University of Pittsburgh)

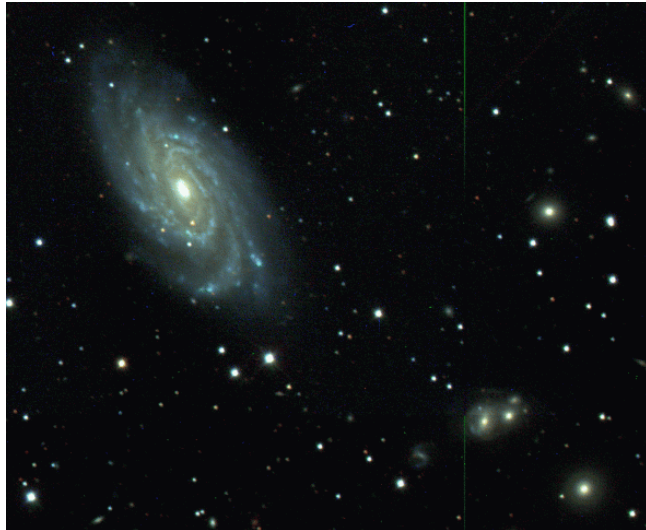
Drug Screening



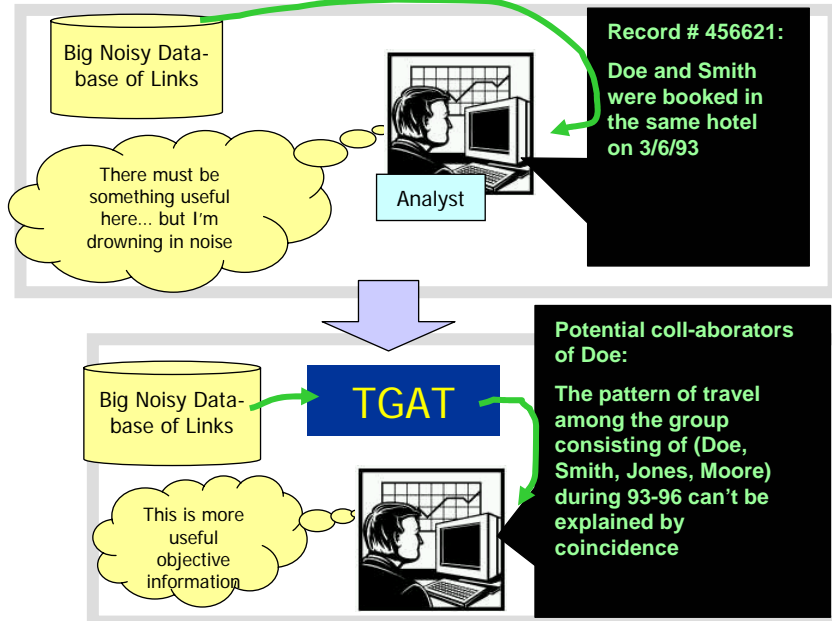
Autonomous self-tweaking engines



Big Astrophysics Automated Science



Intelligence Data



Outline



Cached Sufficient Statistics

Kd-trees and Ball Trees

K-nearest neighbor with ball trees

- Very fast non-parametric classification

- skewed binary outputs

- General binary outputs

- multi-classed outputs

Very fast kernel-based statistics

- n-point computations

- clustering

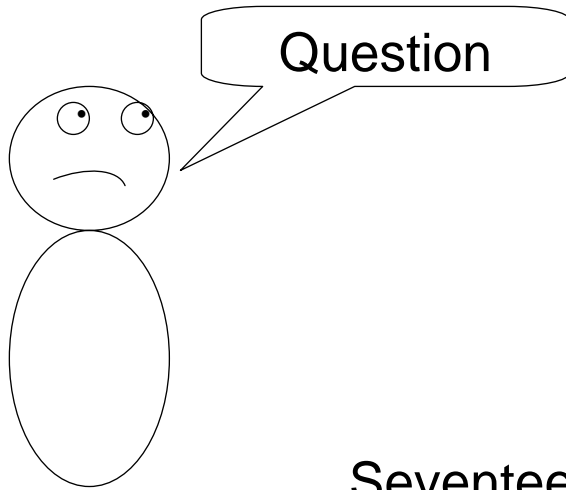
- non-parametric clustering (overdensity hunting)

Active learning for anomaly hunting

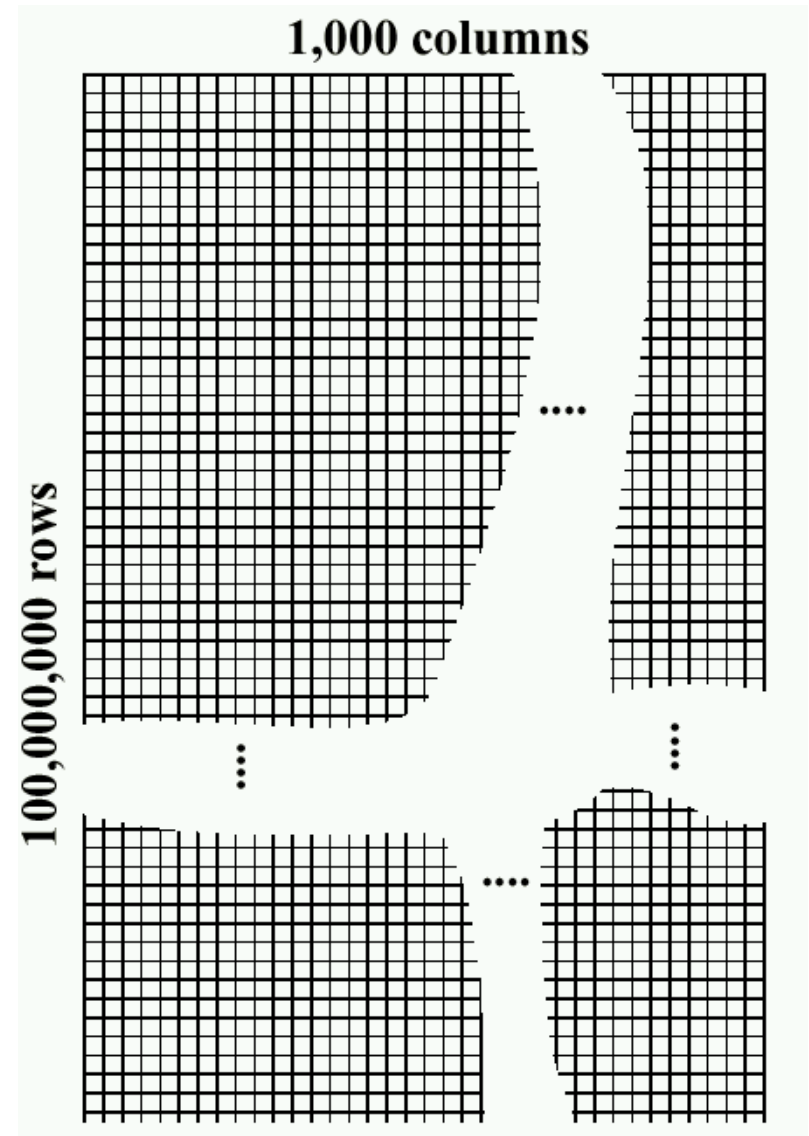
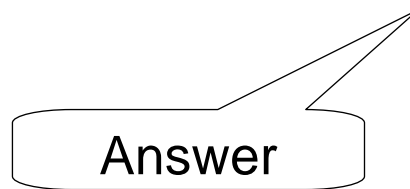
GMorph: Efficient Galaxy morphology fitting

Other Auton topics

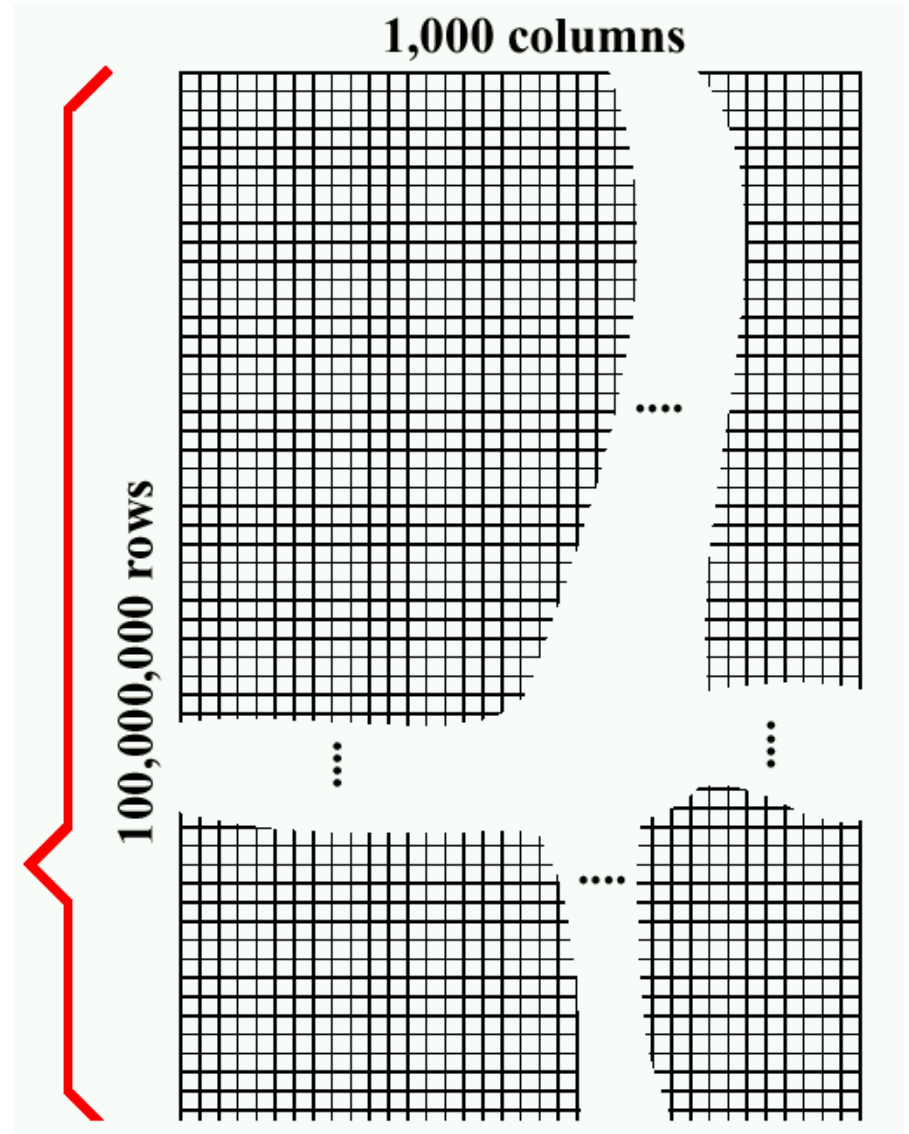
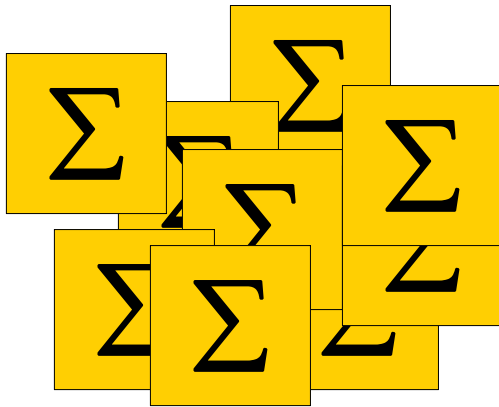
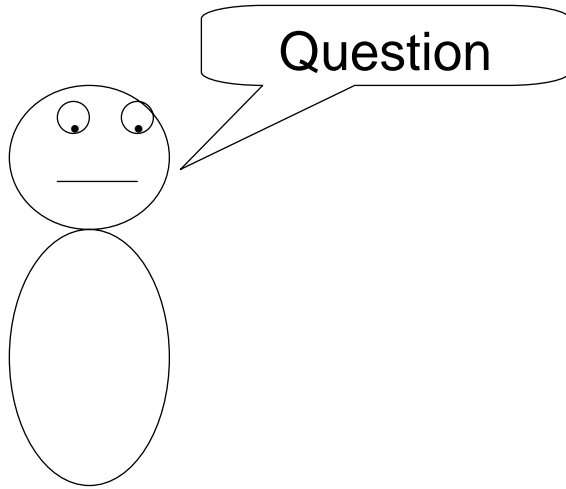
Data Analysis: The new days



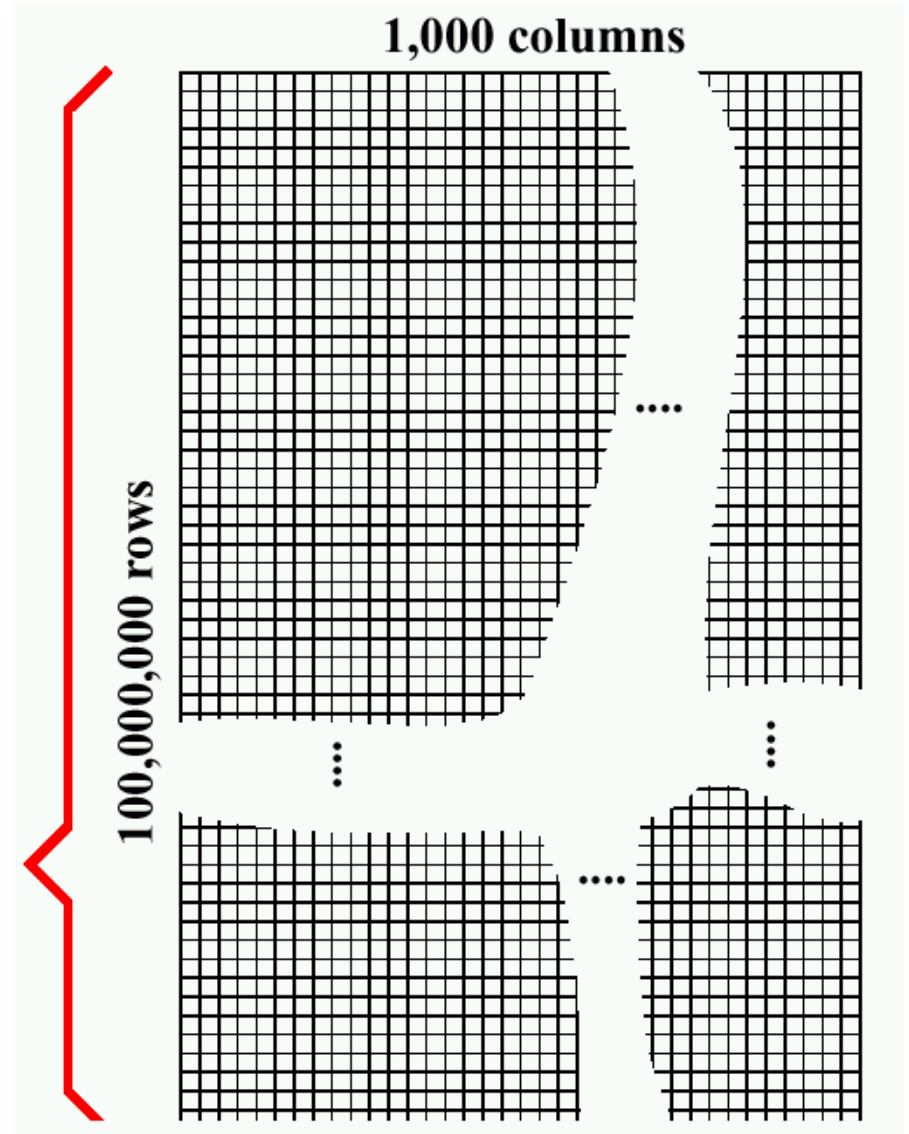
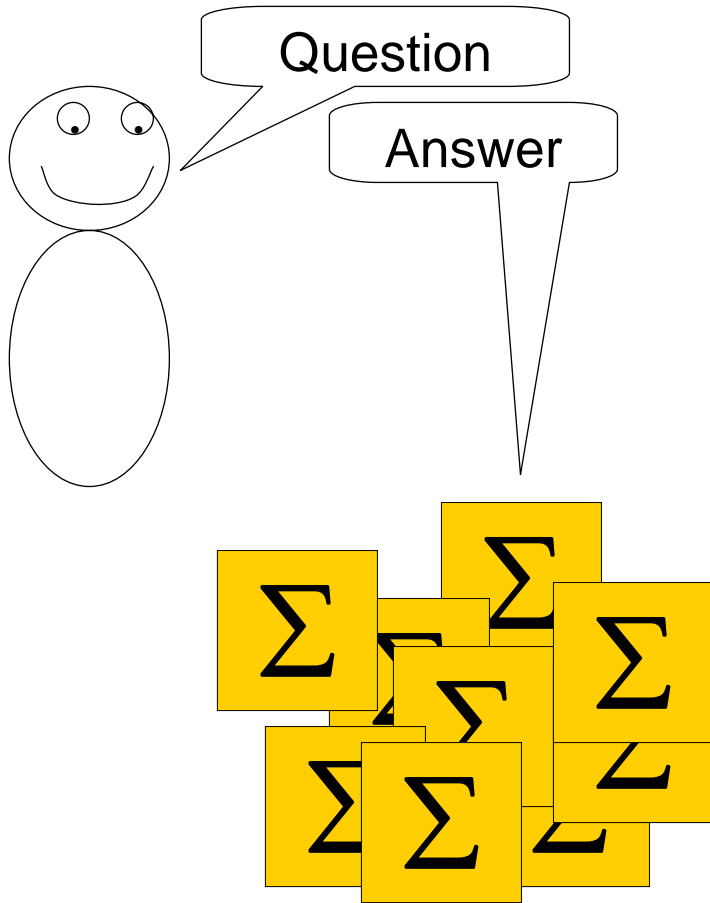
Seventeen months later...



Cached Sufficient Statistics

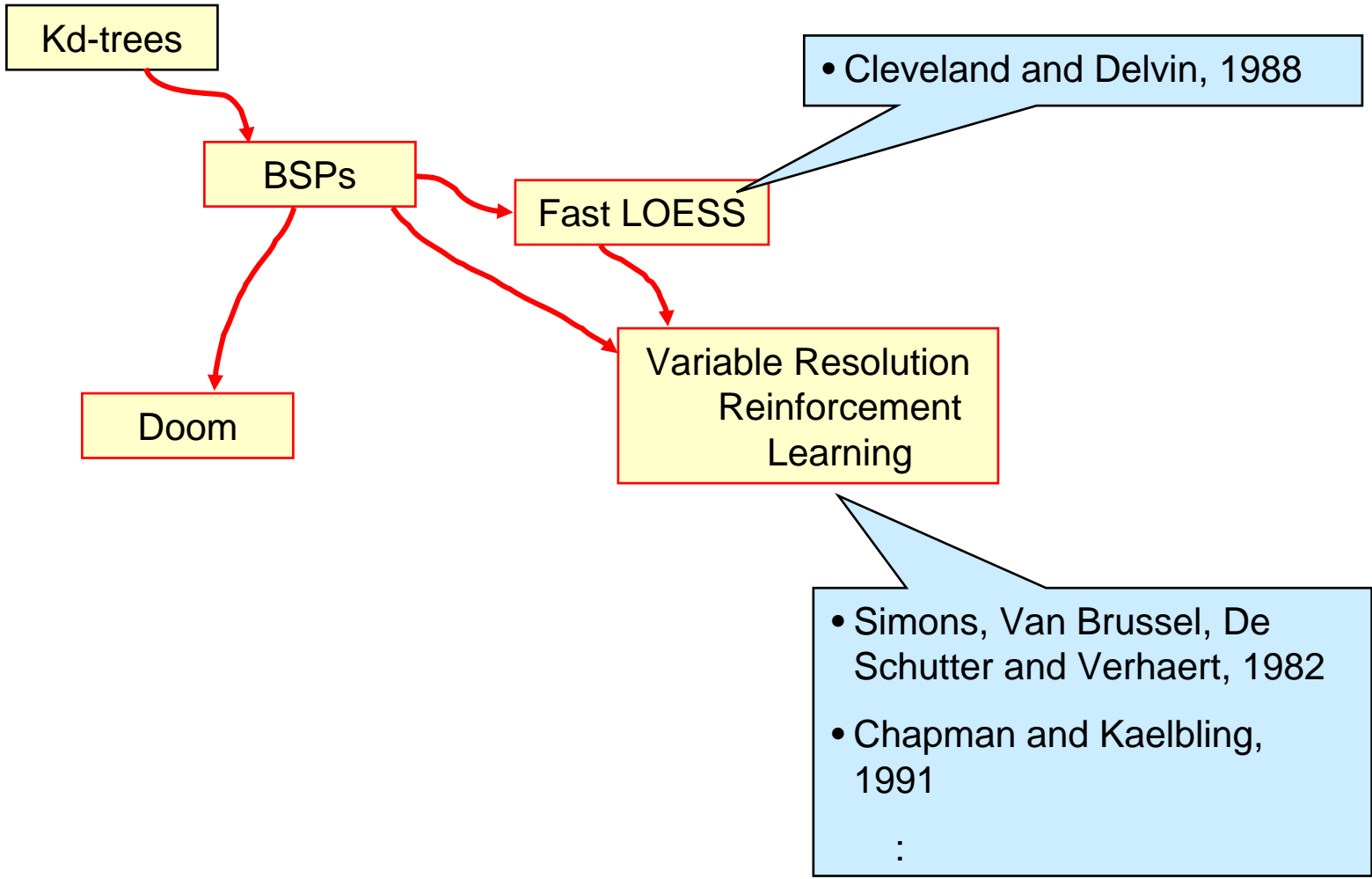


Cached Sufficient Statistics



Kd-trees

- Friedman, Bentley and Finkel, 1977



Kd-trees

BSPs

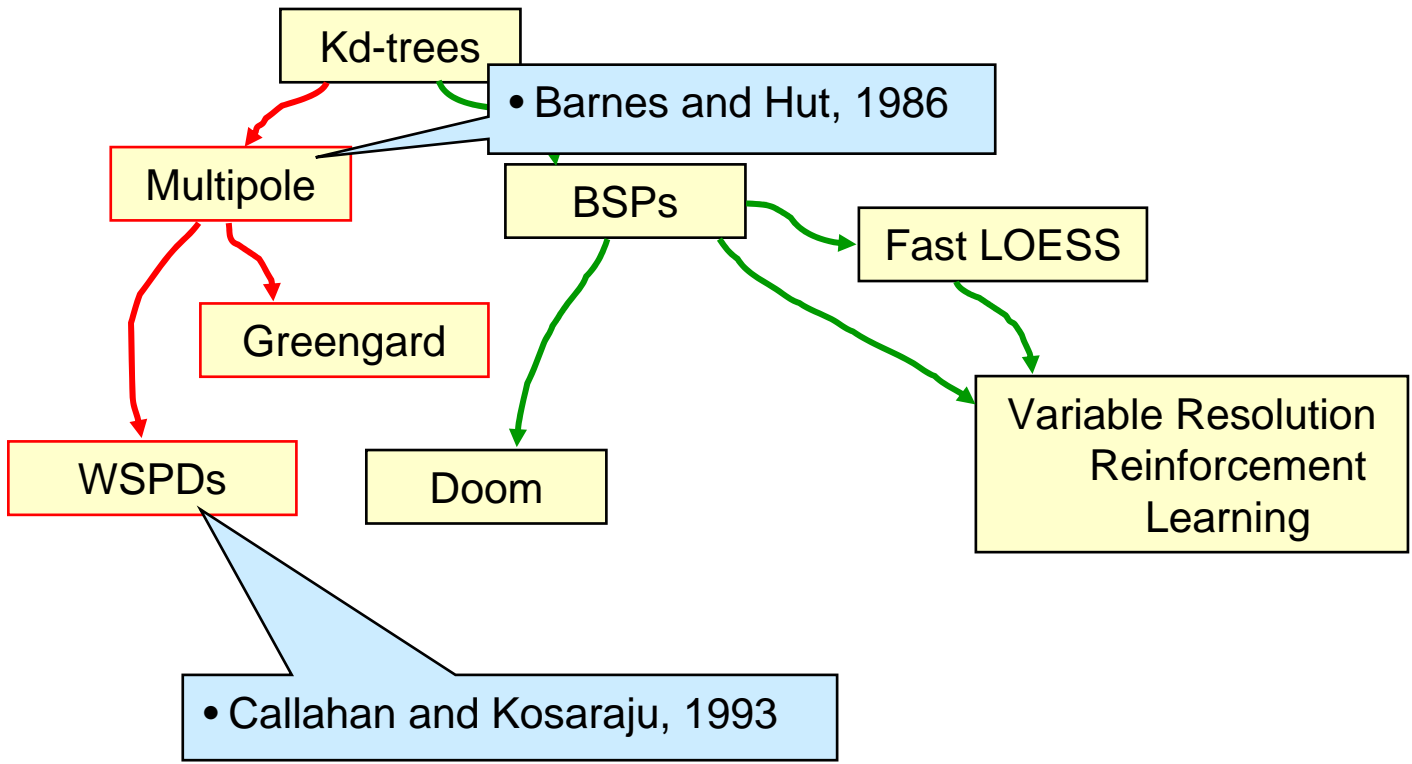
Fast LOESS

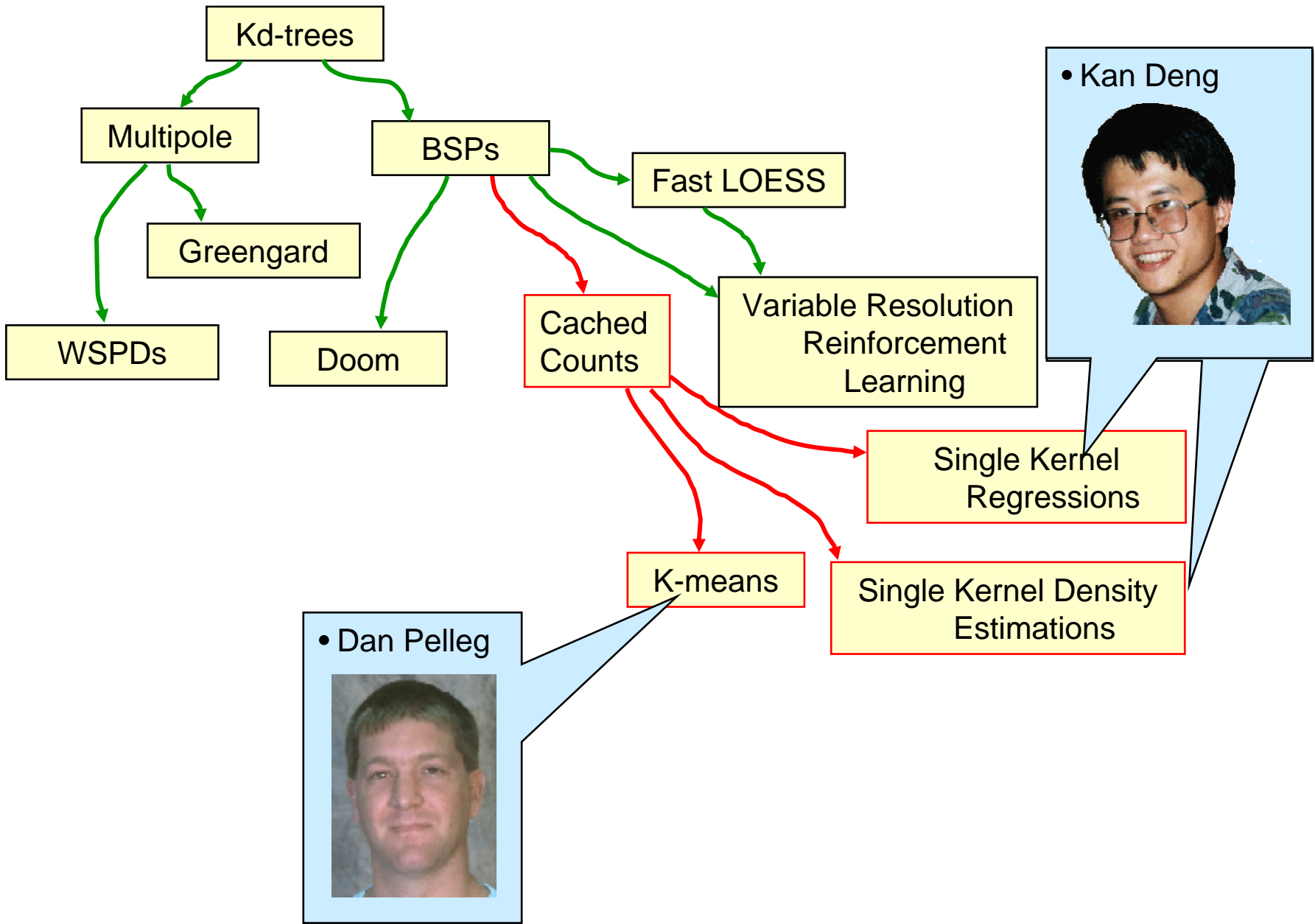
Doom

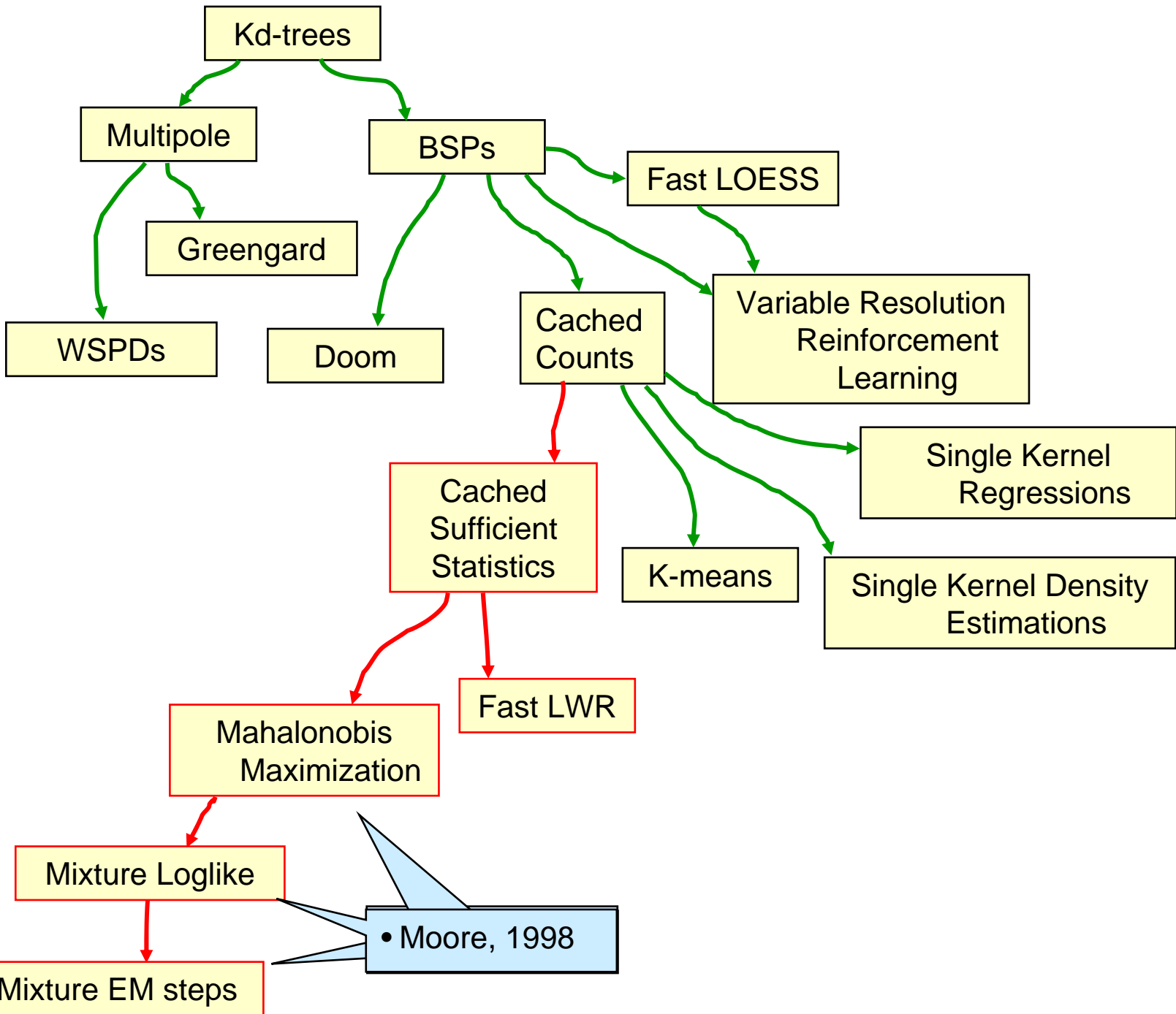
Variable Resolution
Reinforcement
Learning

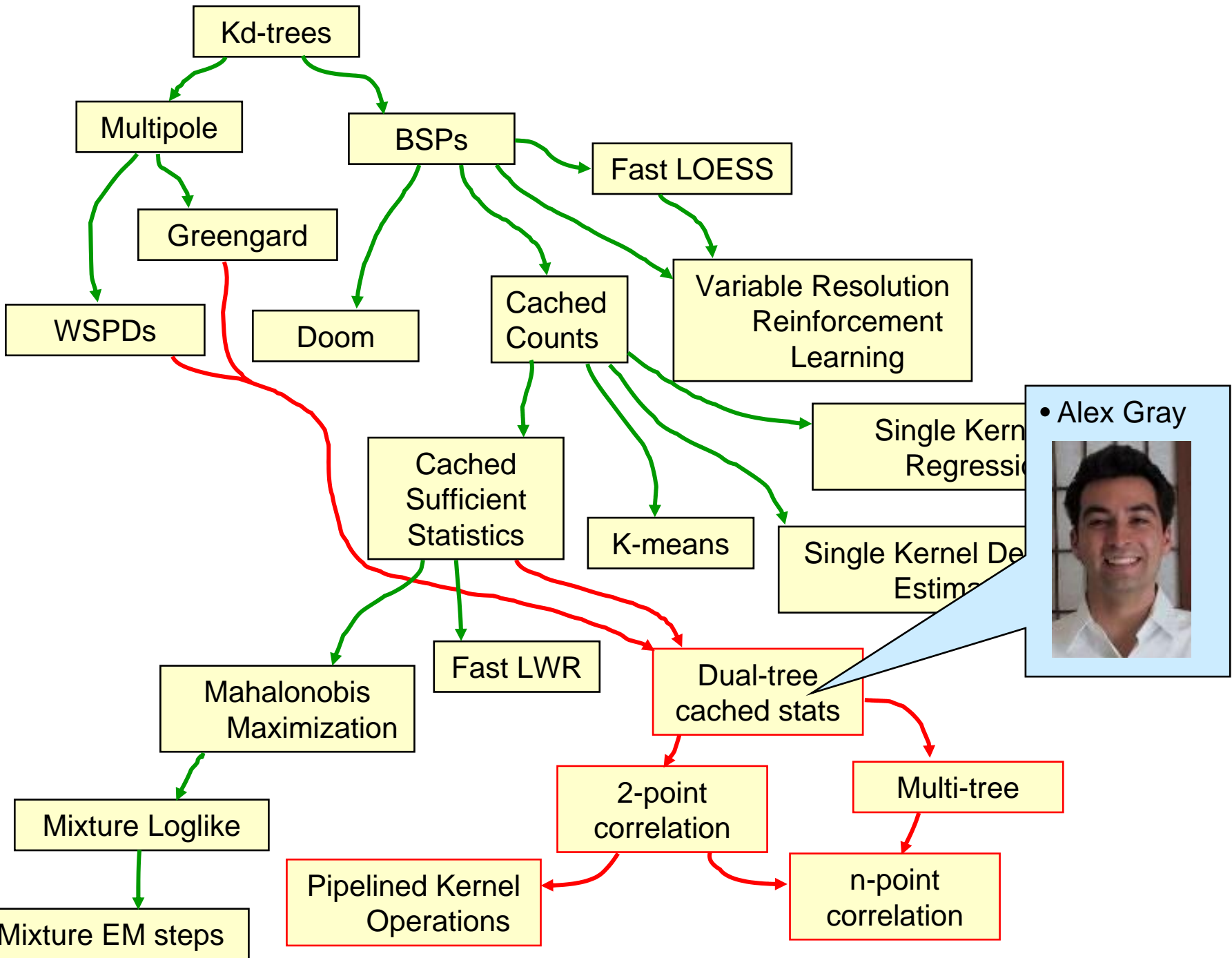
• Cleveland and Delvin, 1988

• Simons, Van Brussel, De Schutter and Verhaert, 1982
• Chapman and Kaelbling, 1991
:



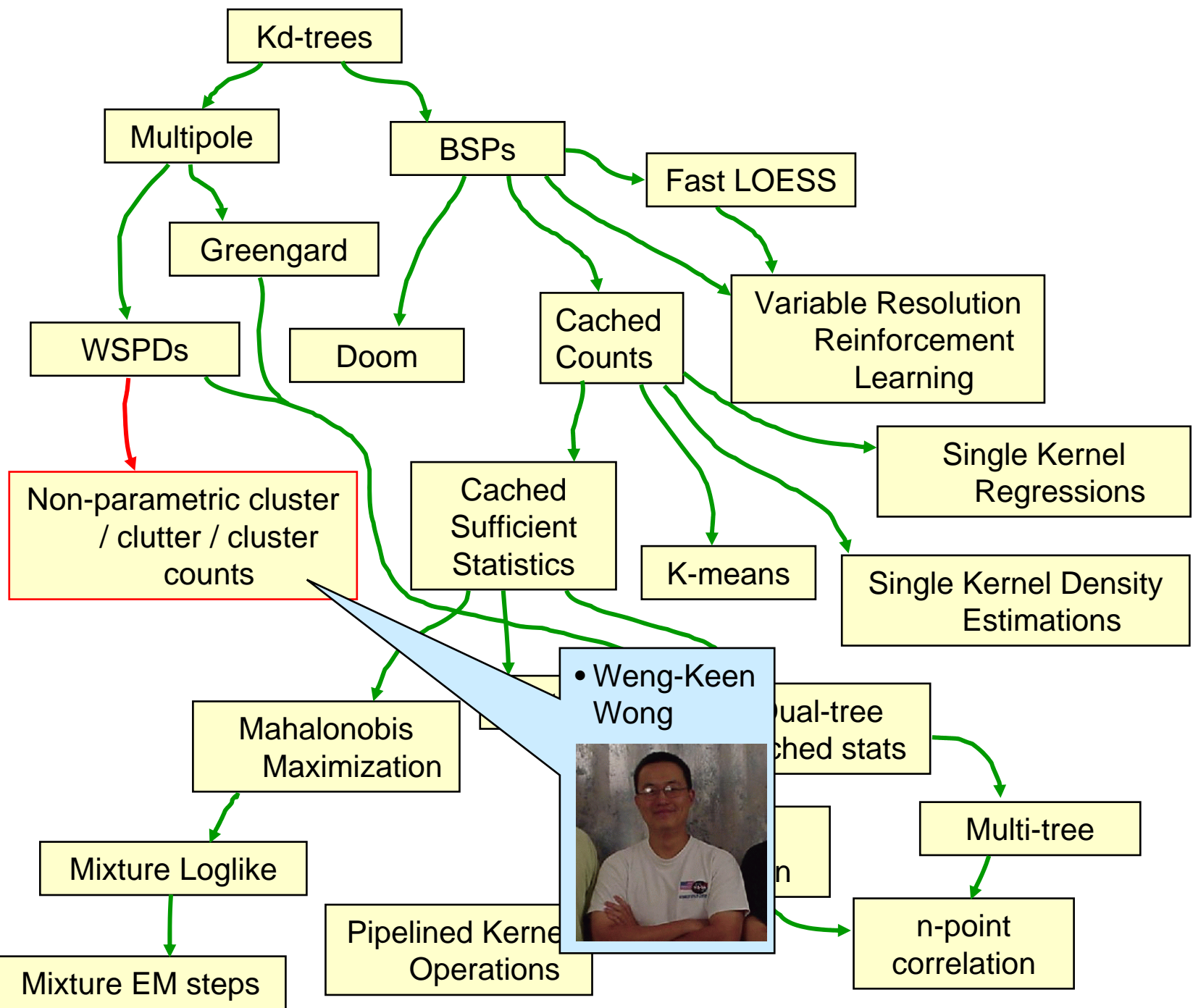






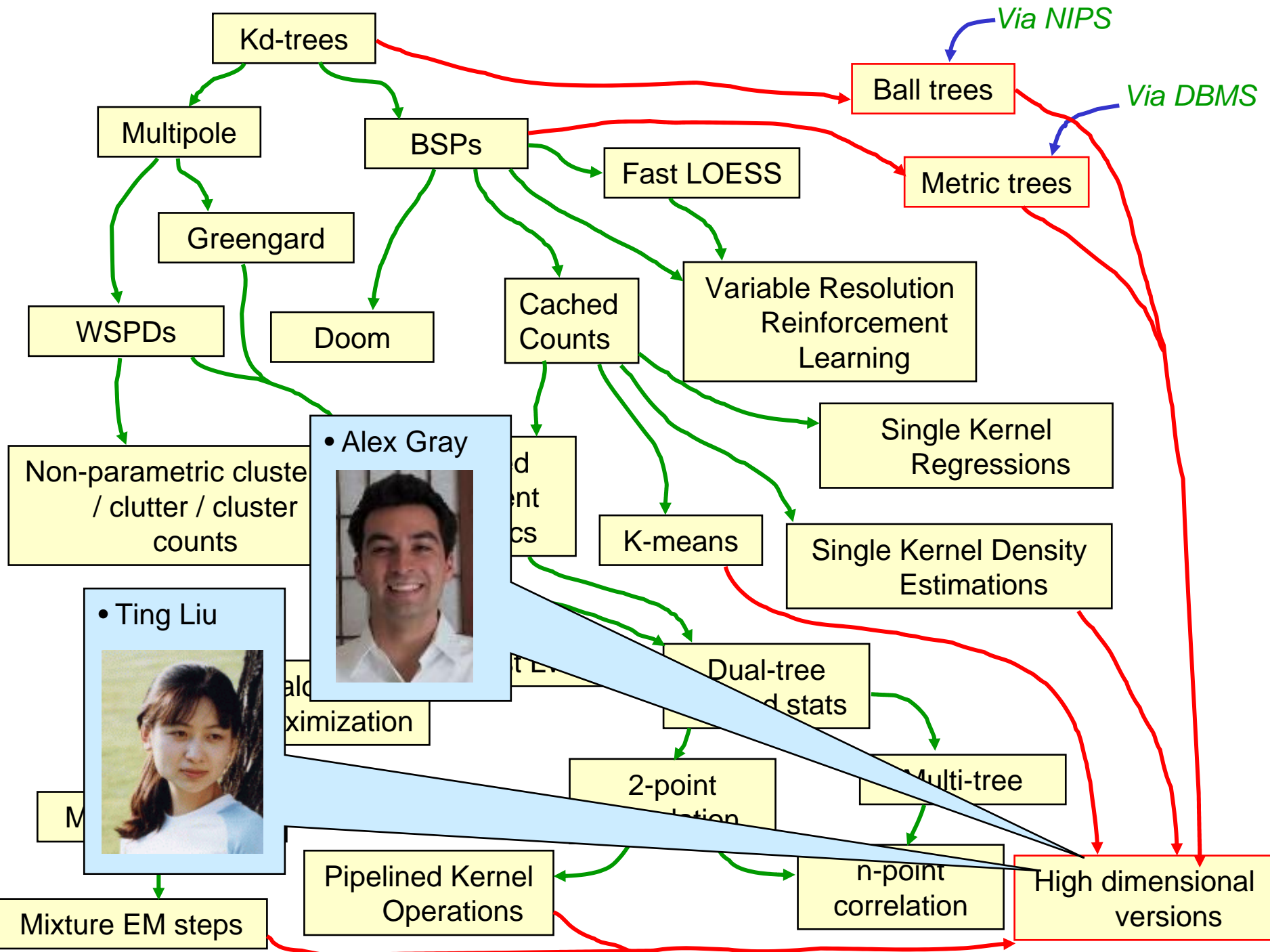
• Alex Gray





• Weng-Keen Wong





Outline

Cached Sufficient Statistics

▶ Ball Trees (= Metric Trees)

K-nearest neighbor with ball trees

Very fast non-parametric classification

skewed binary outputs

General binary outputs

multi-classed outputs

Very fast kernel-based statistics

n-point computations

clustering

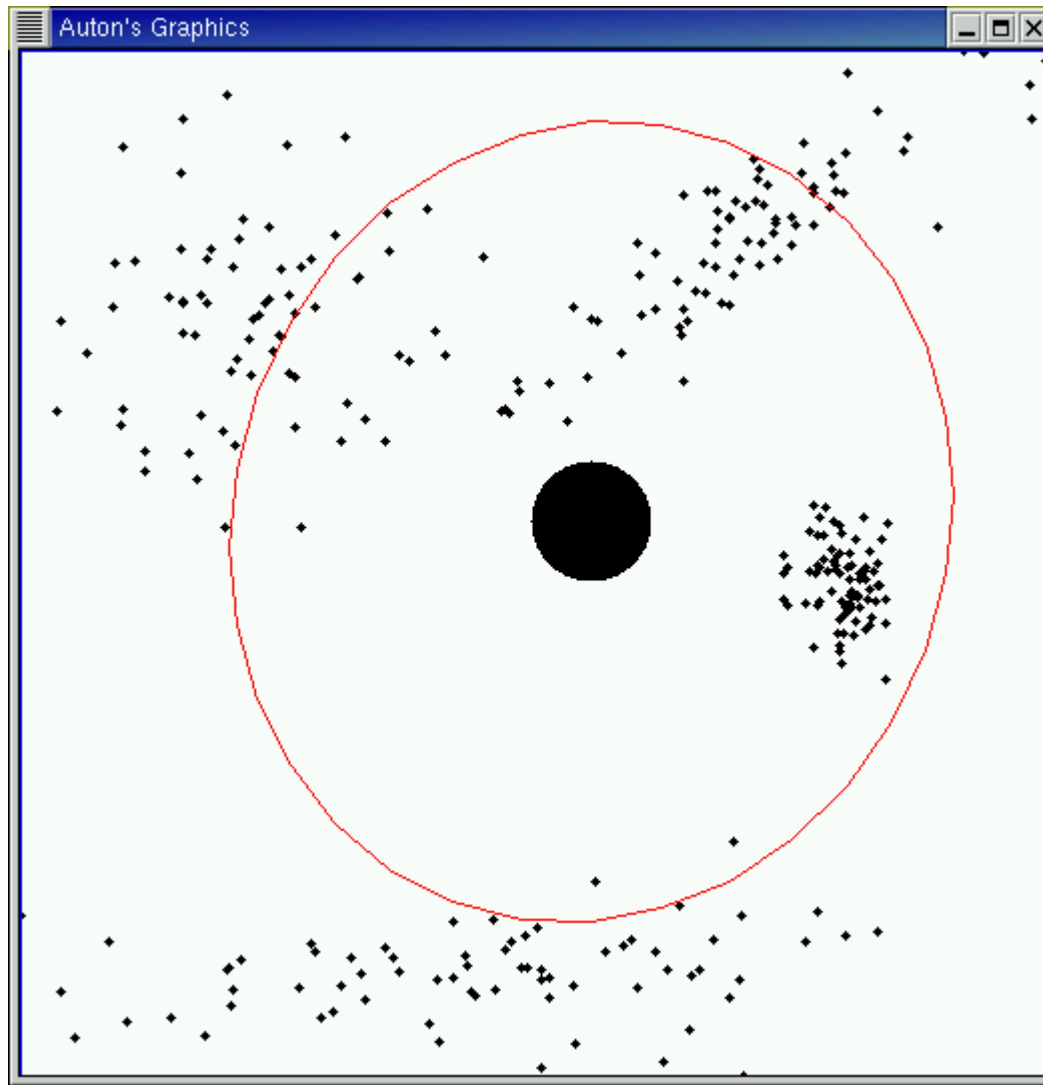
non-parametric clustering (overdensity hunting)

Active learning for anomaly hunting

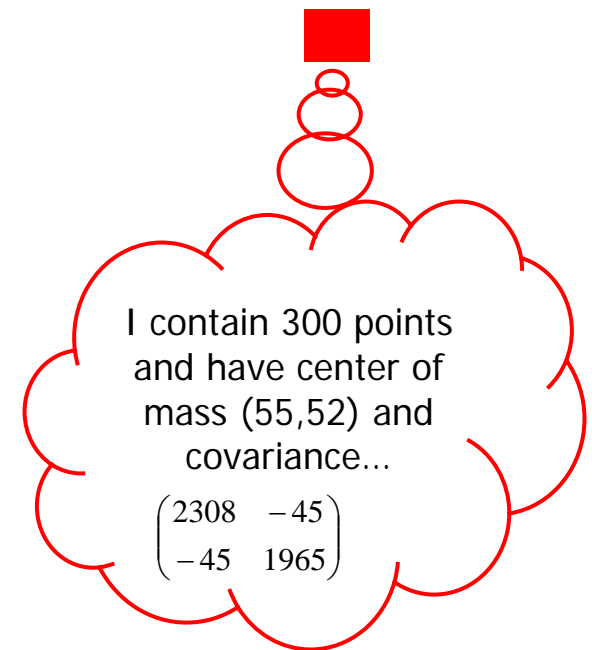
GMorph: Efficient Galaxy morphology fitting

Other Auton topics

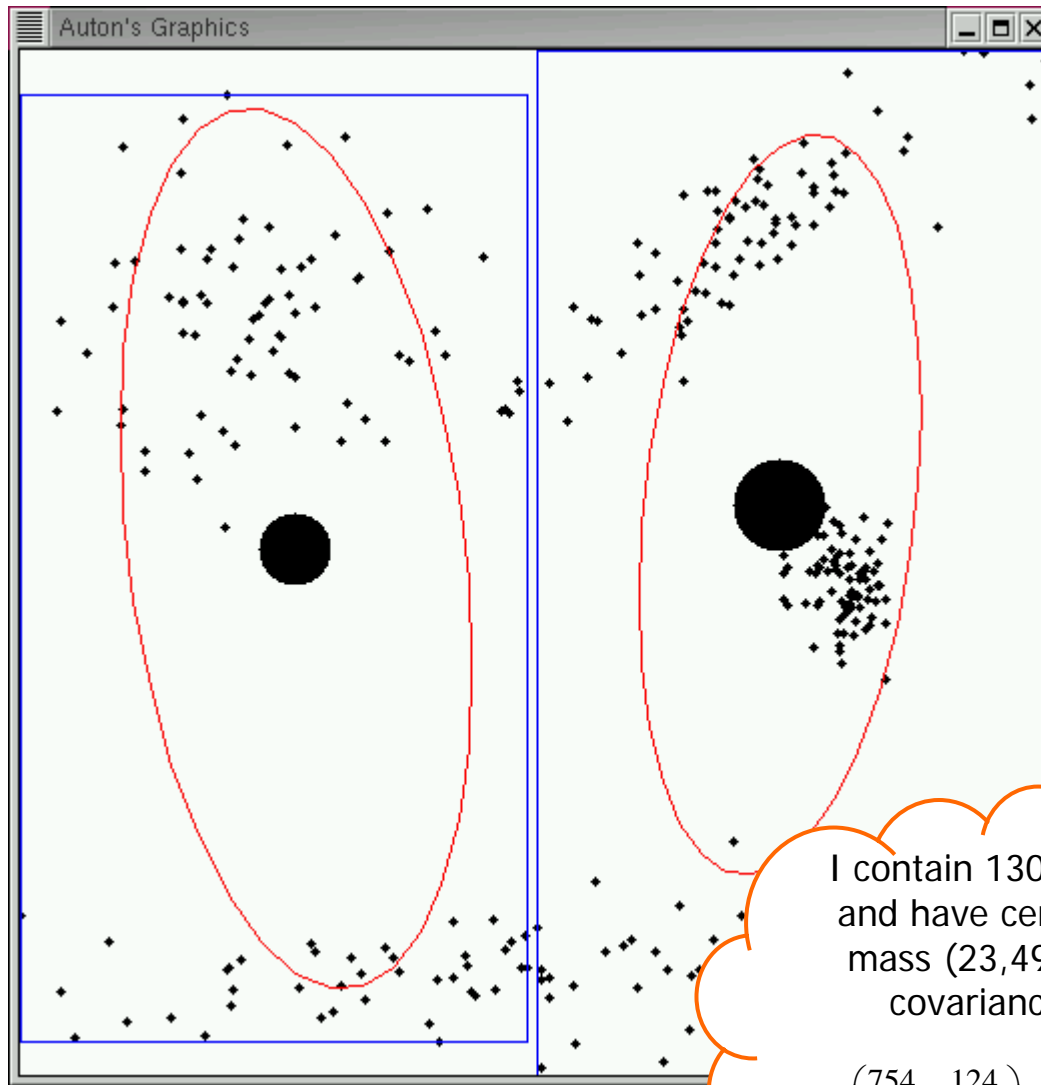
Structure of kdtrees



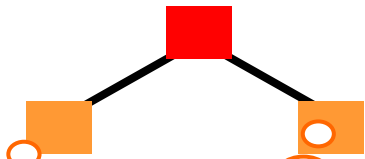
In memory...



Structure of of kdtrees



In memory...



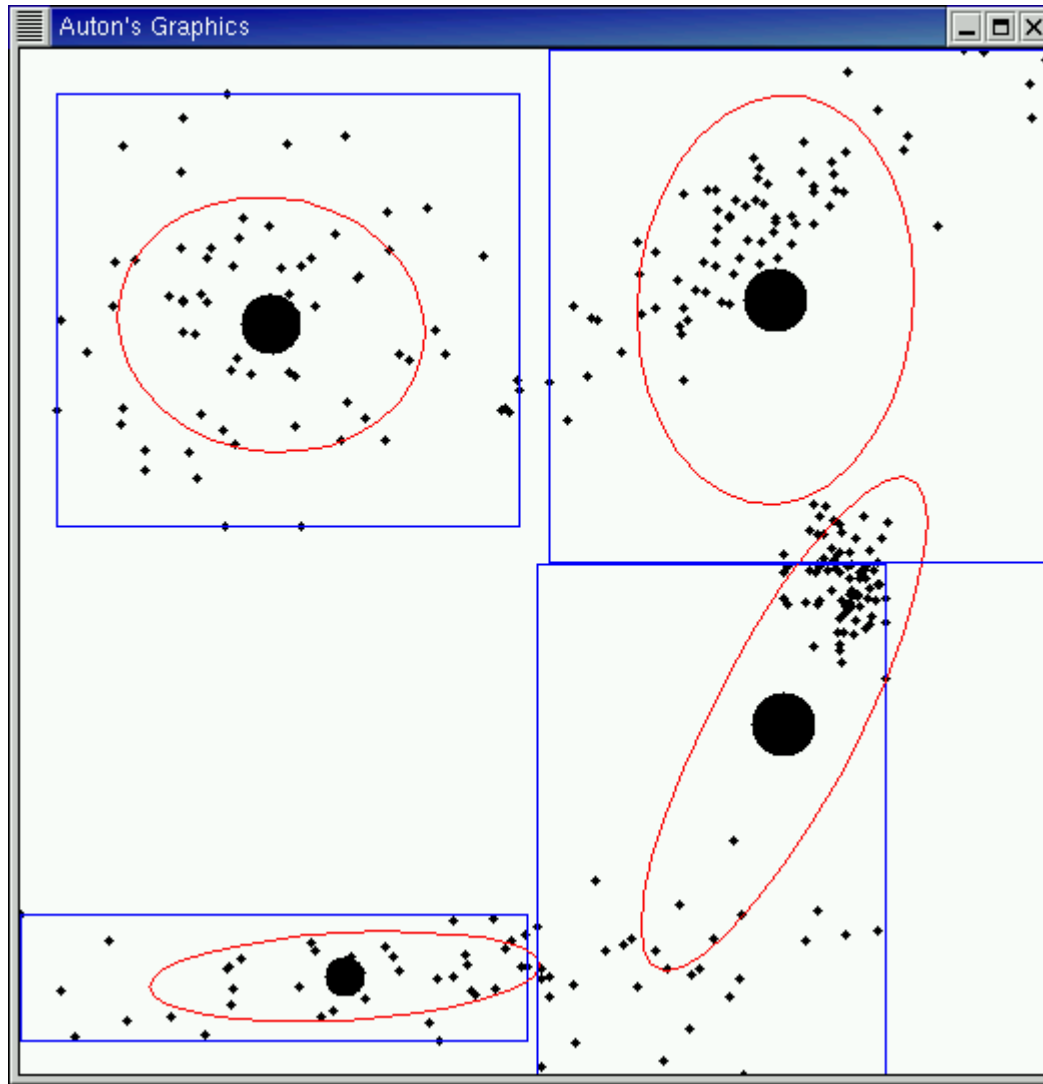
I contain 130 points
and have center of
mass (23,49) and
covariance...

$$\begin{pmatrix} 754 & 124 \\ 124 & 1965 \end{pmatrix}$$

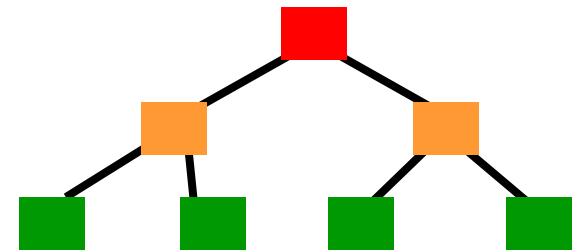
I contain 170 points
and have center of
mass (73,58) and
covariance...

$$\begin{pmatrix} 654 & 24 \\ 24 & 1885 \end{pmatrix}$$

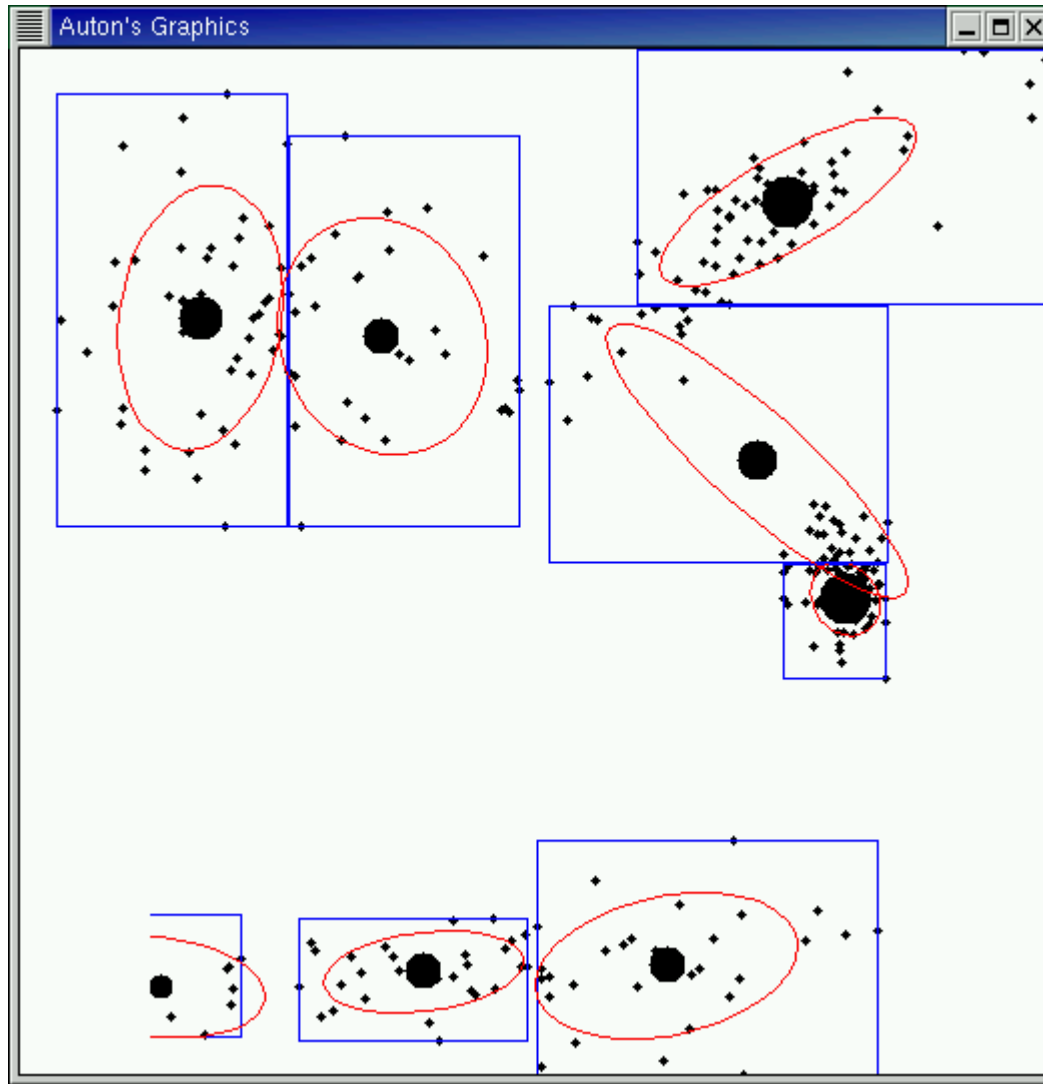
Structure of of kdtrees



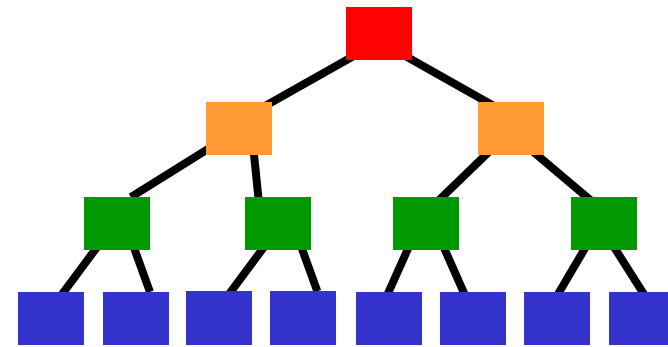
In memory...



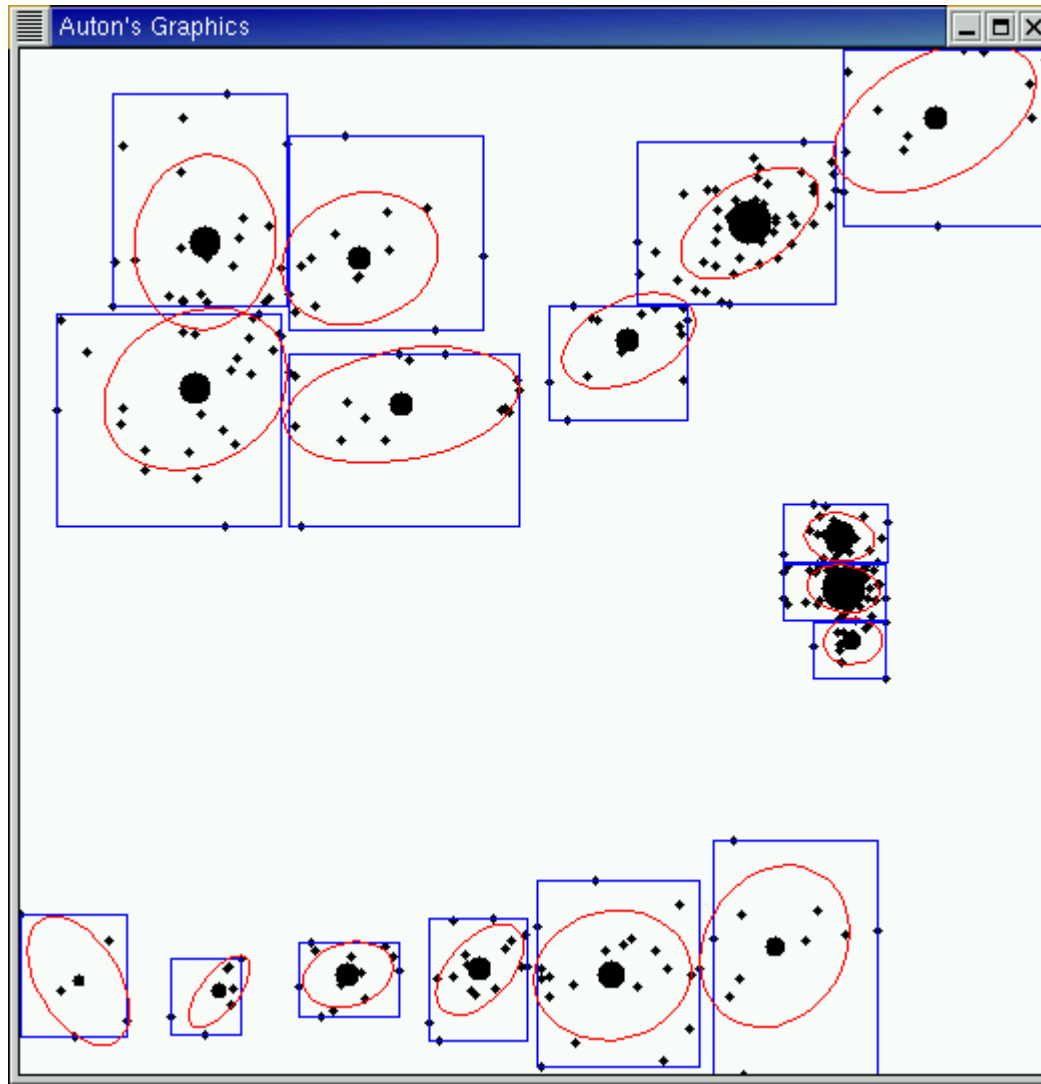
Structure of of kdtrees



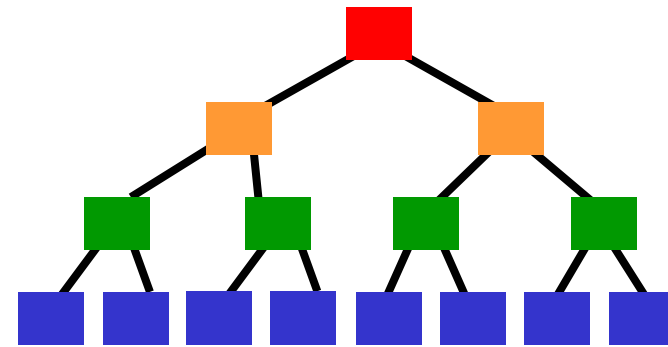
In memory...



Structure of of kdtrees

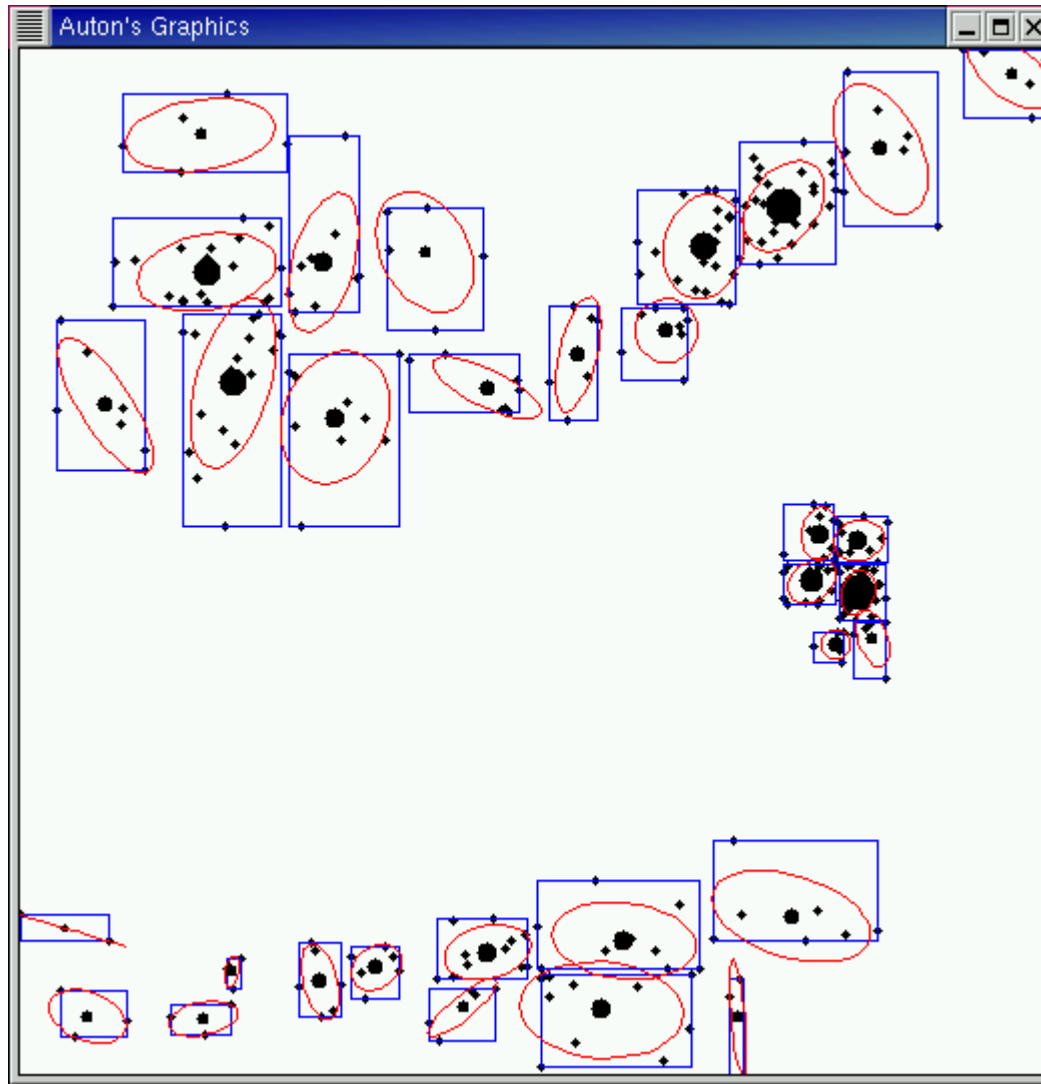


In memory...

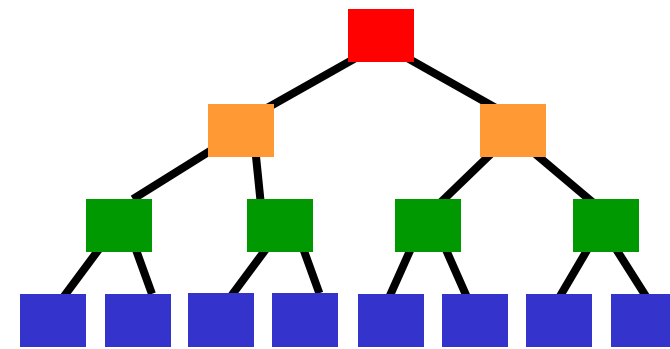


<more levels here>

Structure of of kdtrees

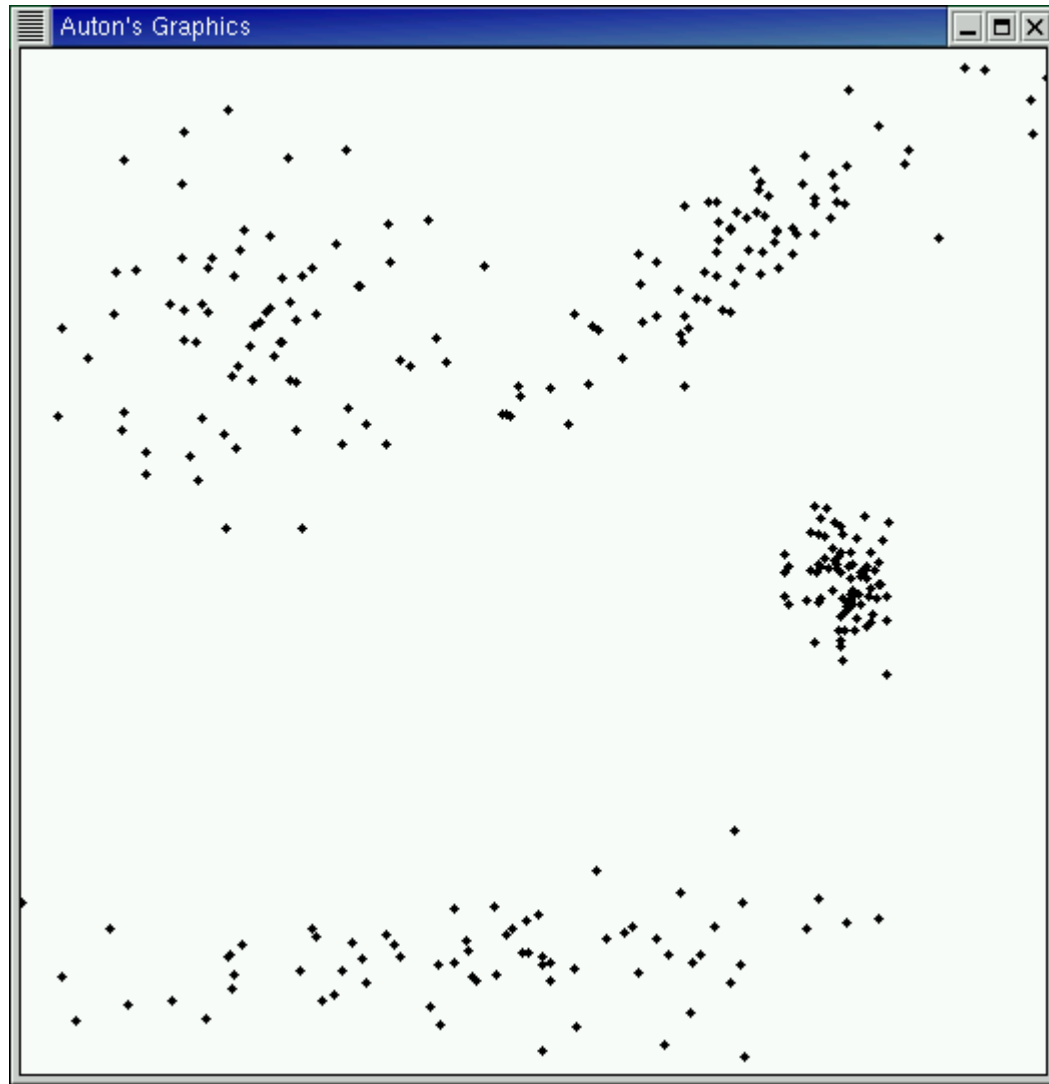


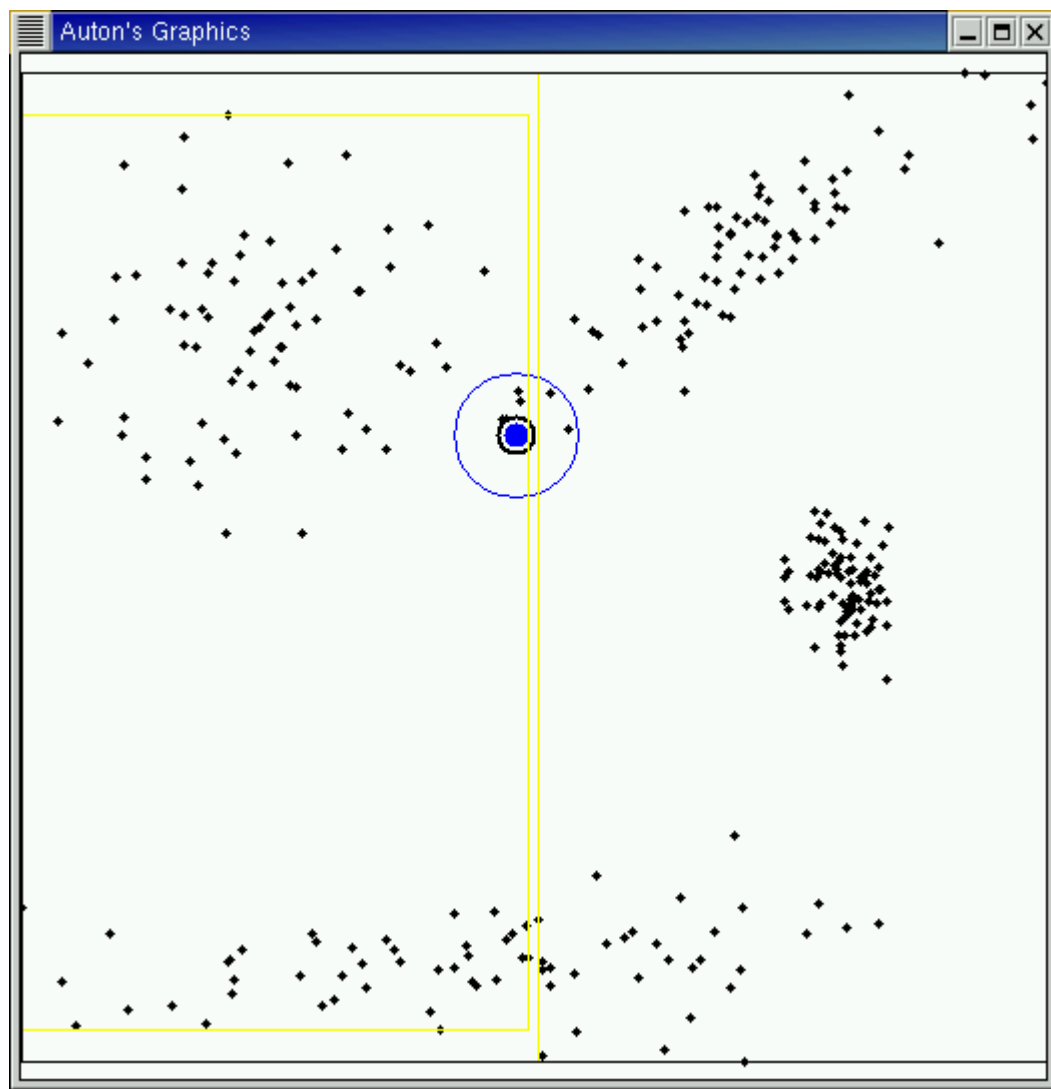
In memory...

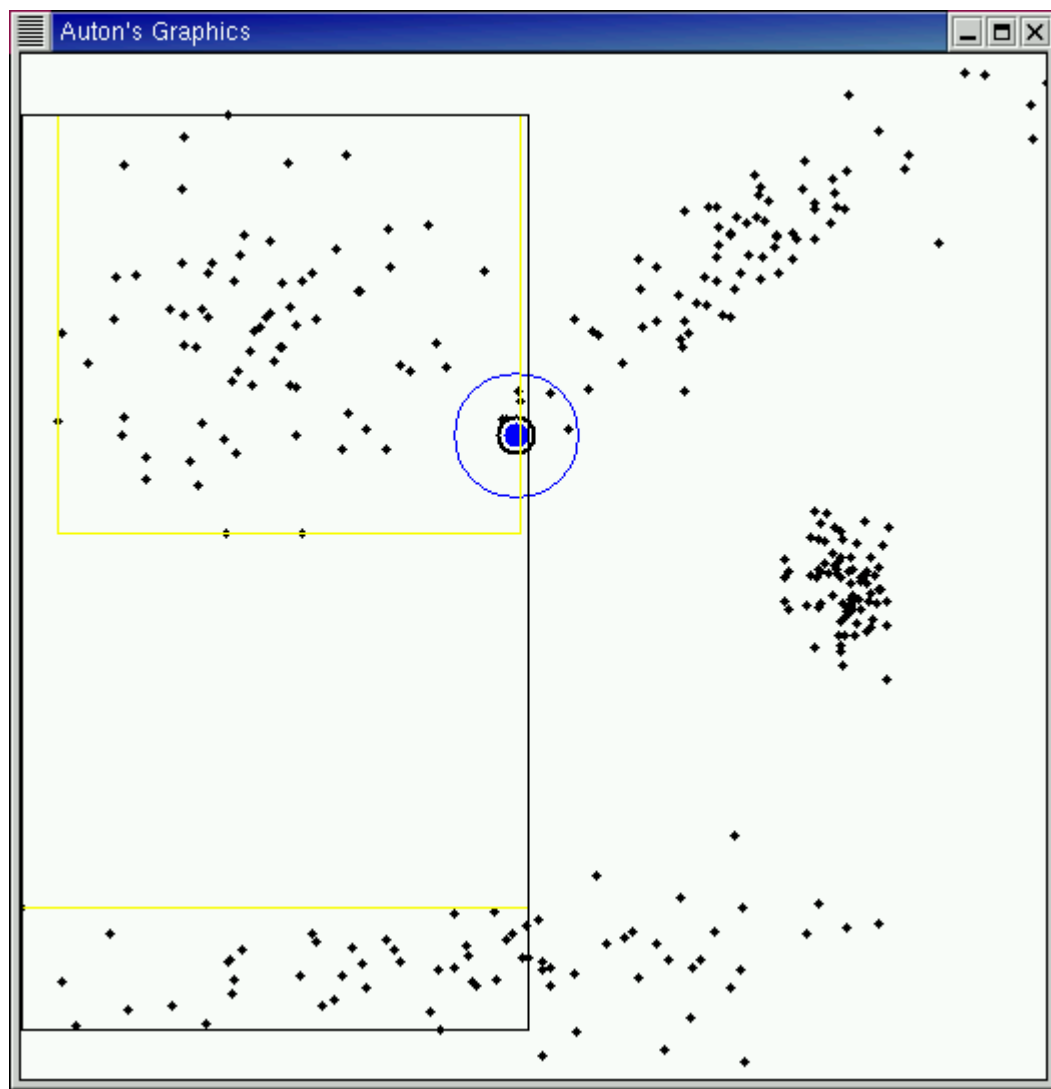


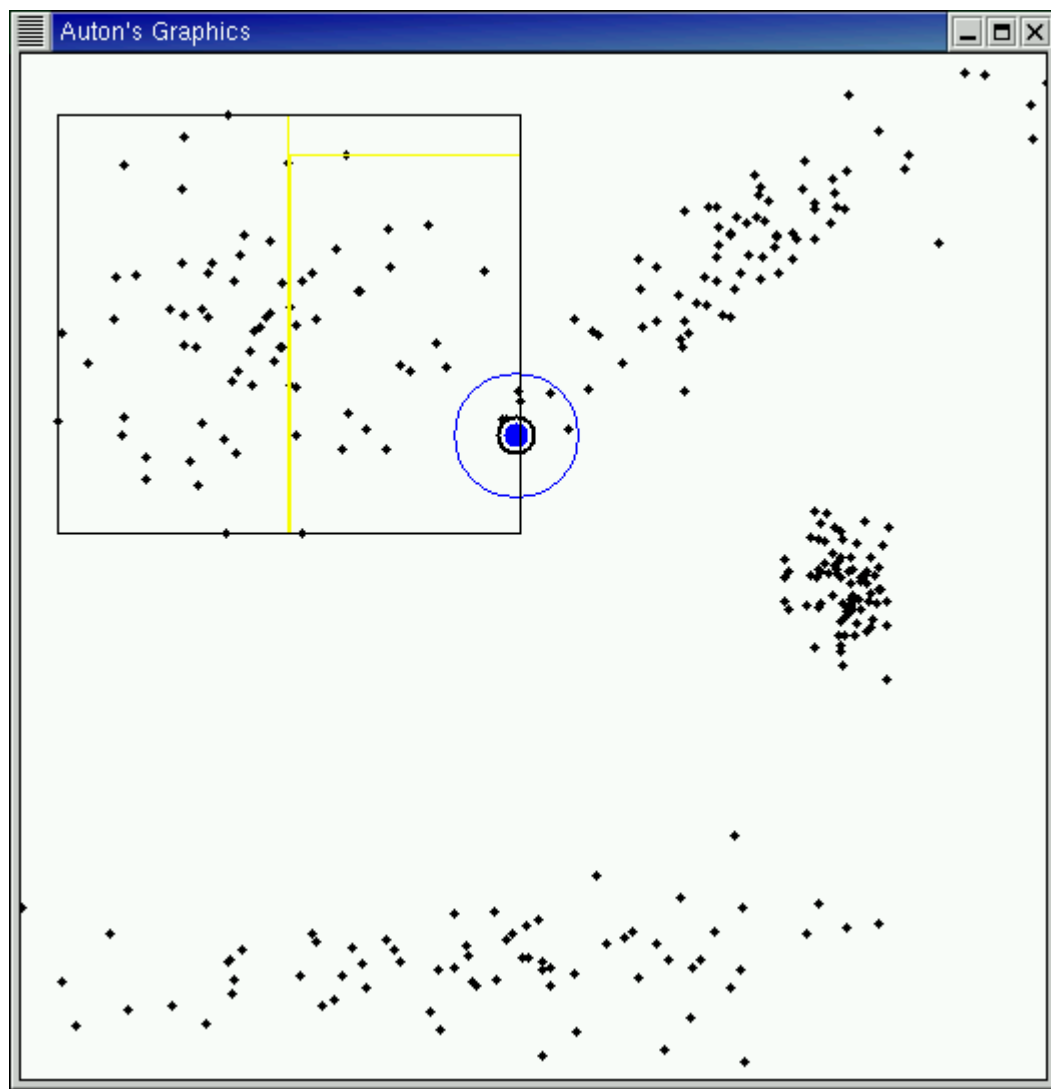
<more levels here>

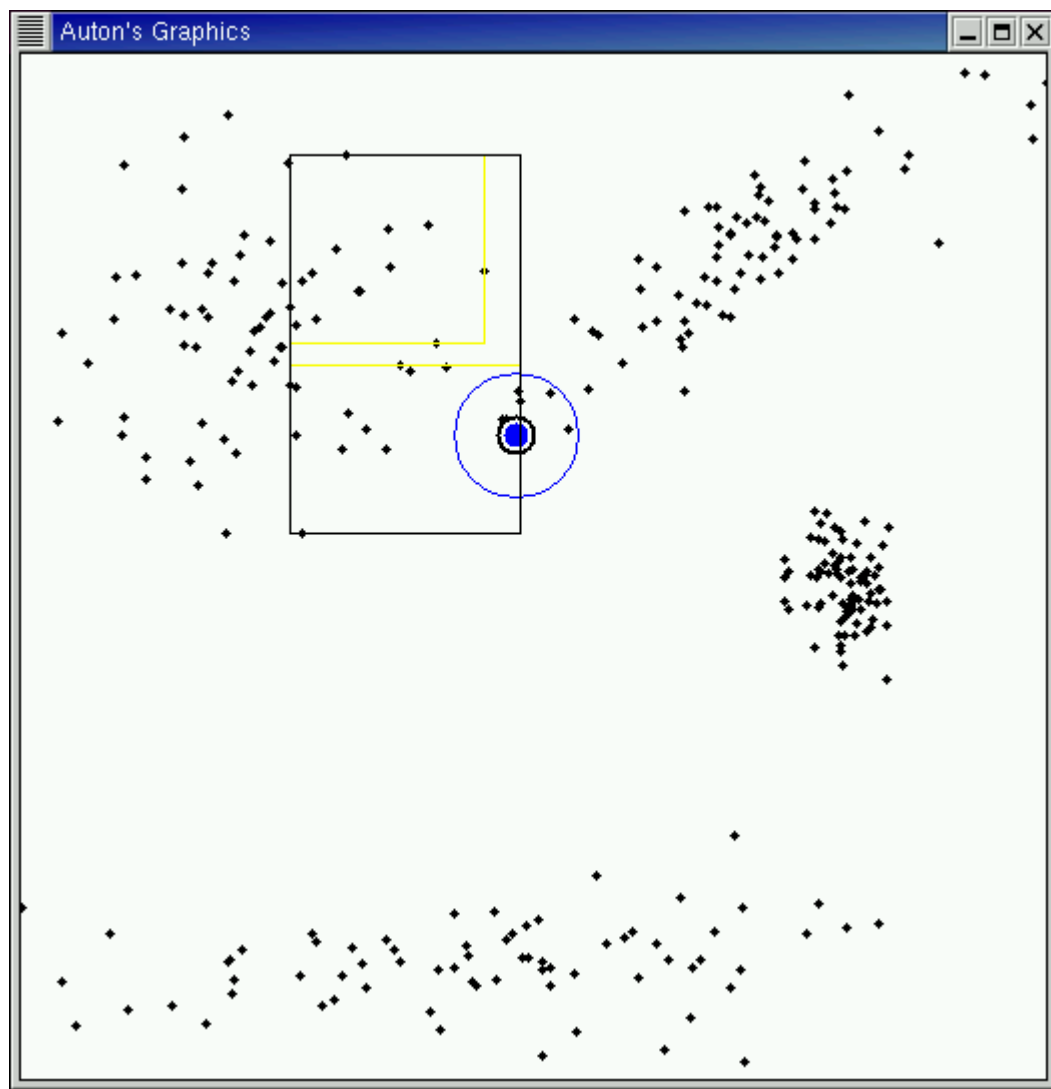
Range Search

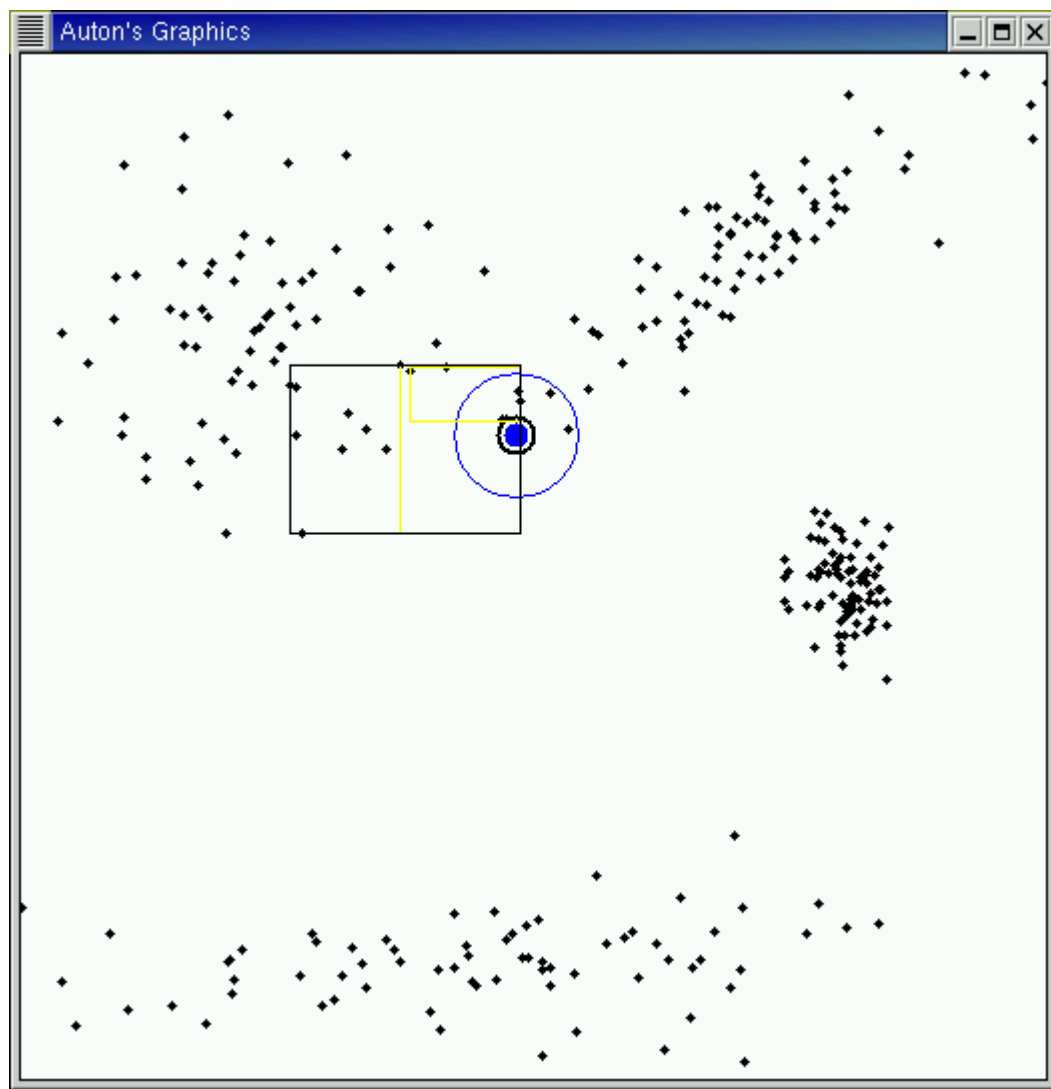


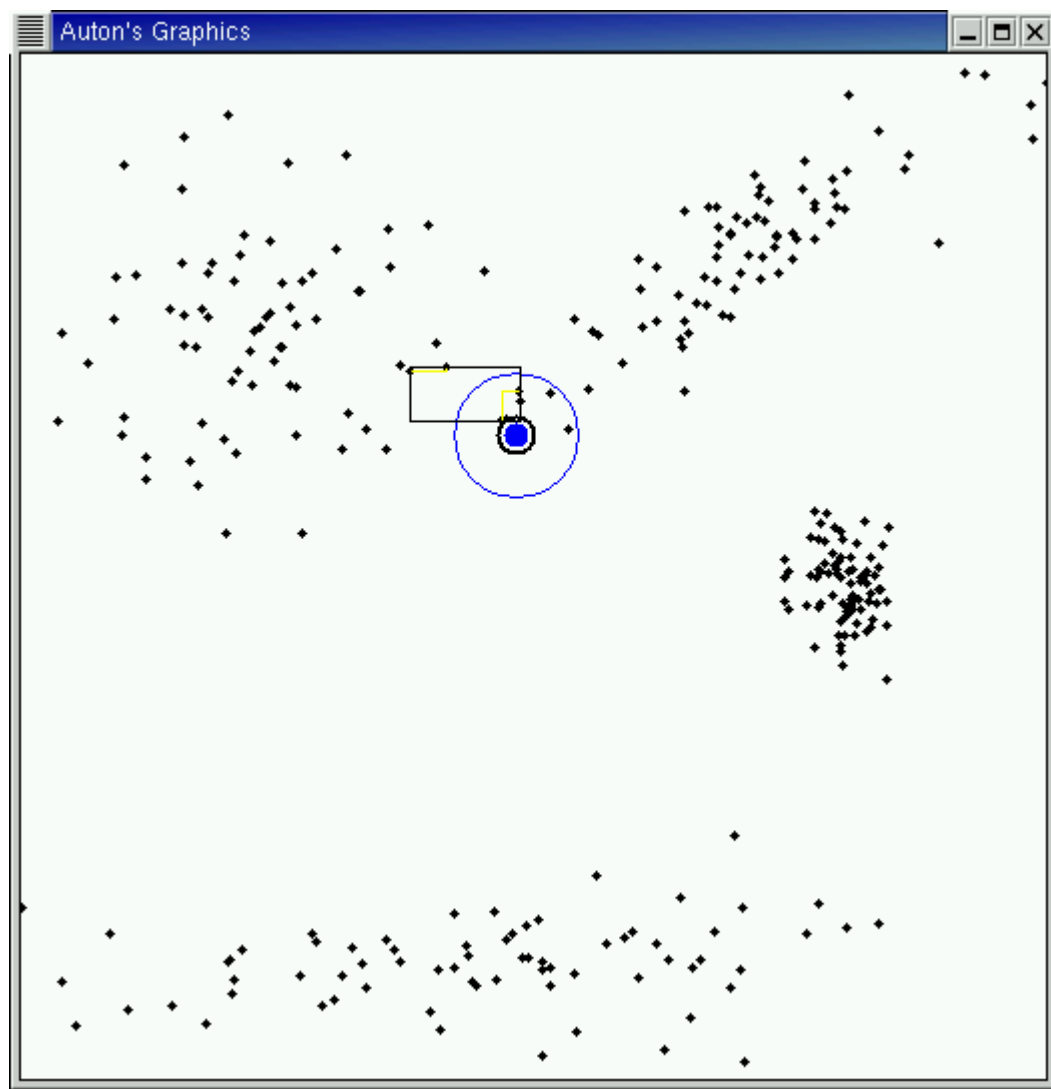


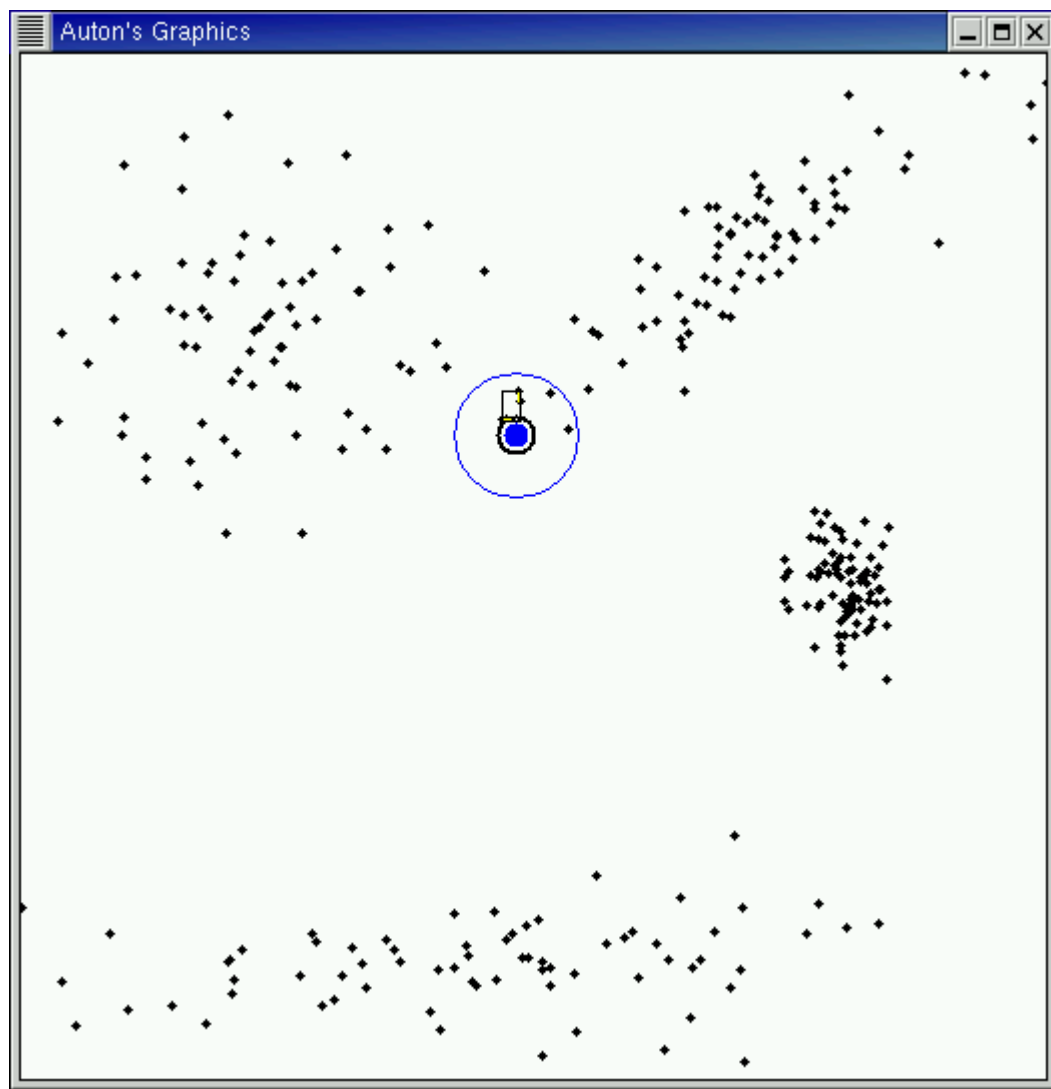


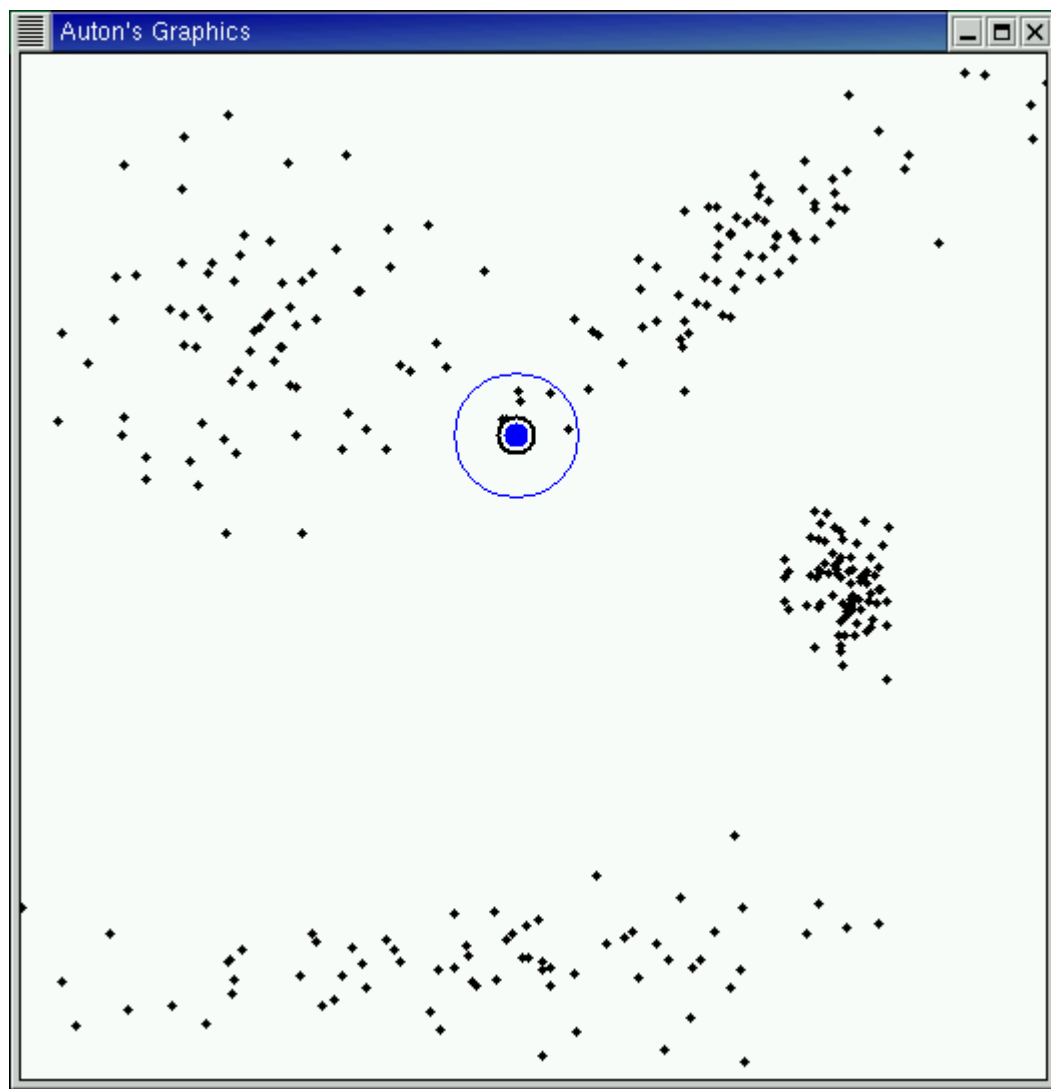


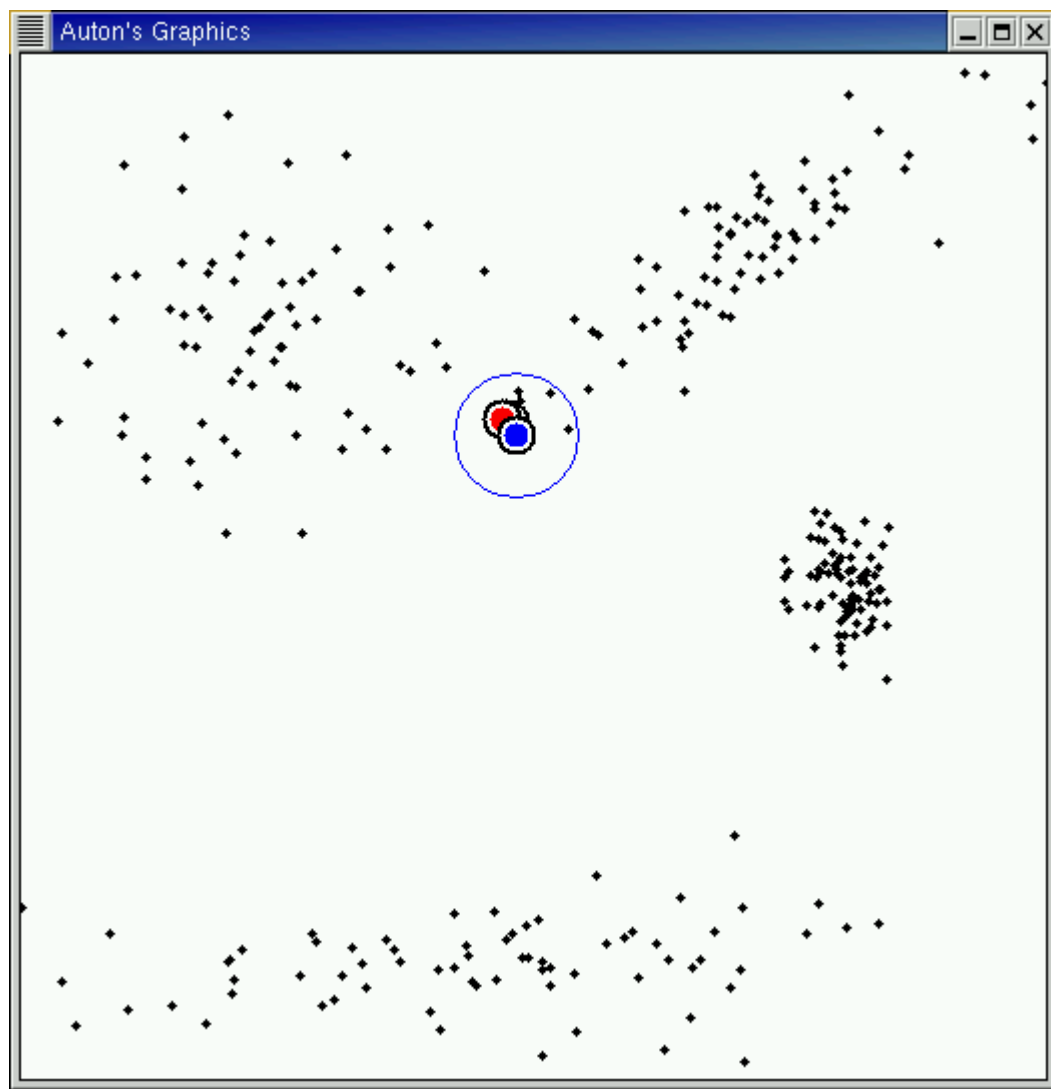


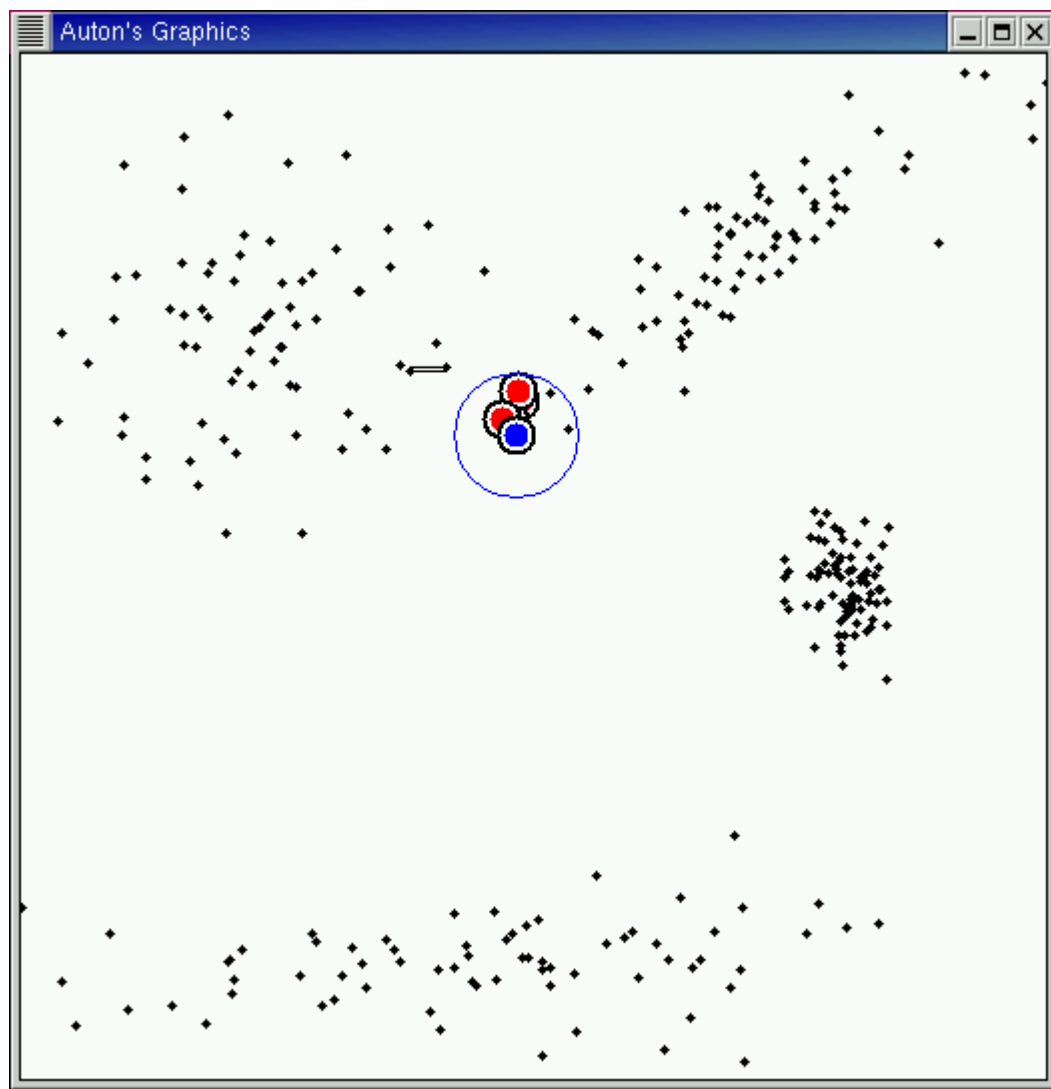


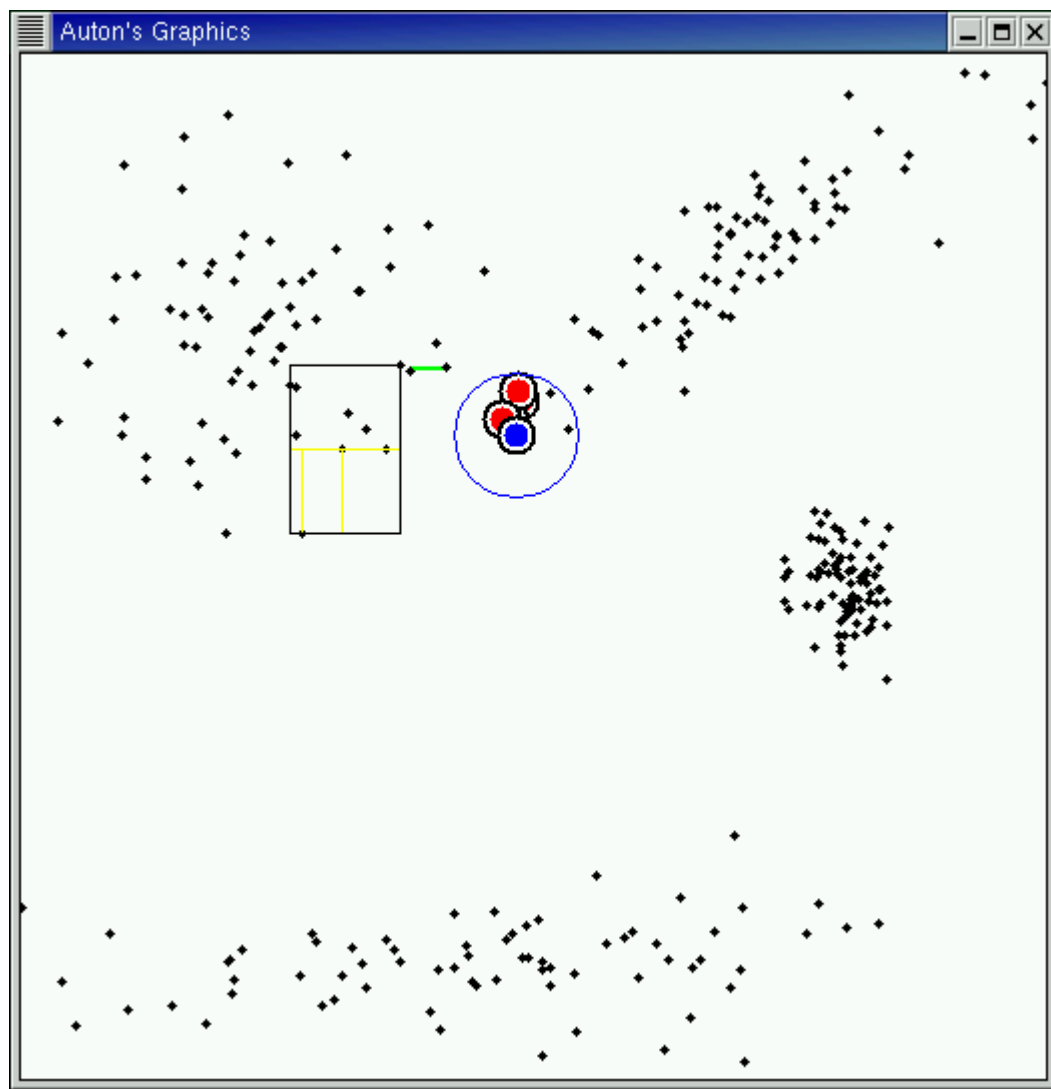


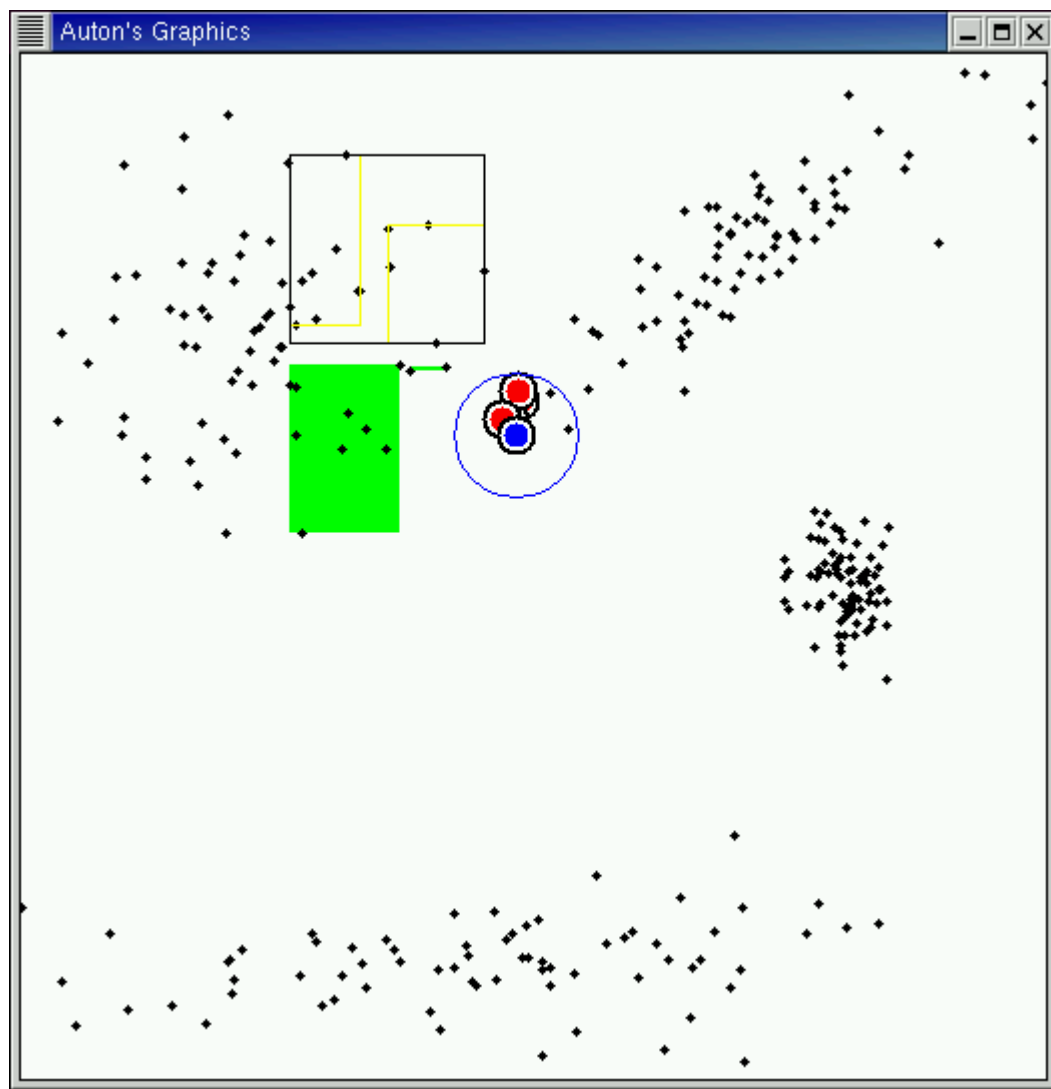


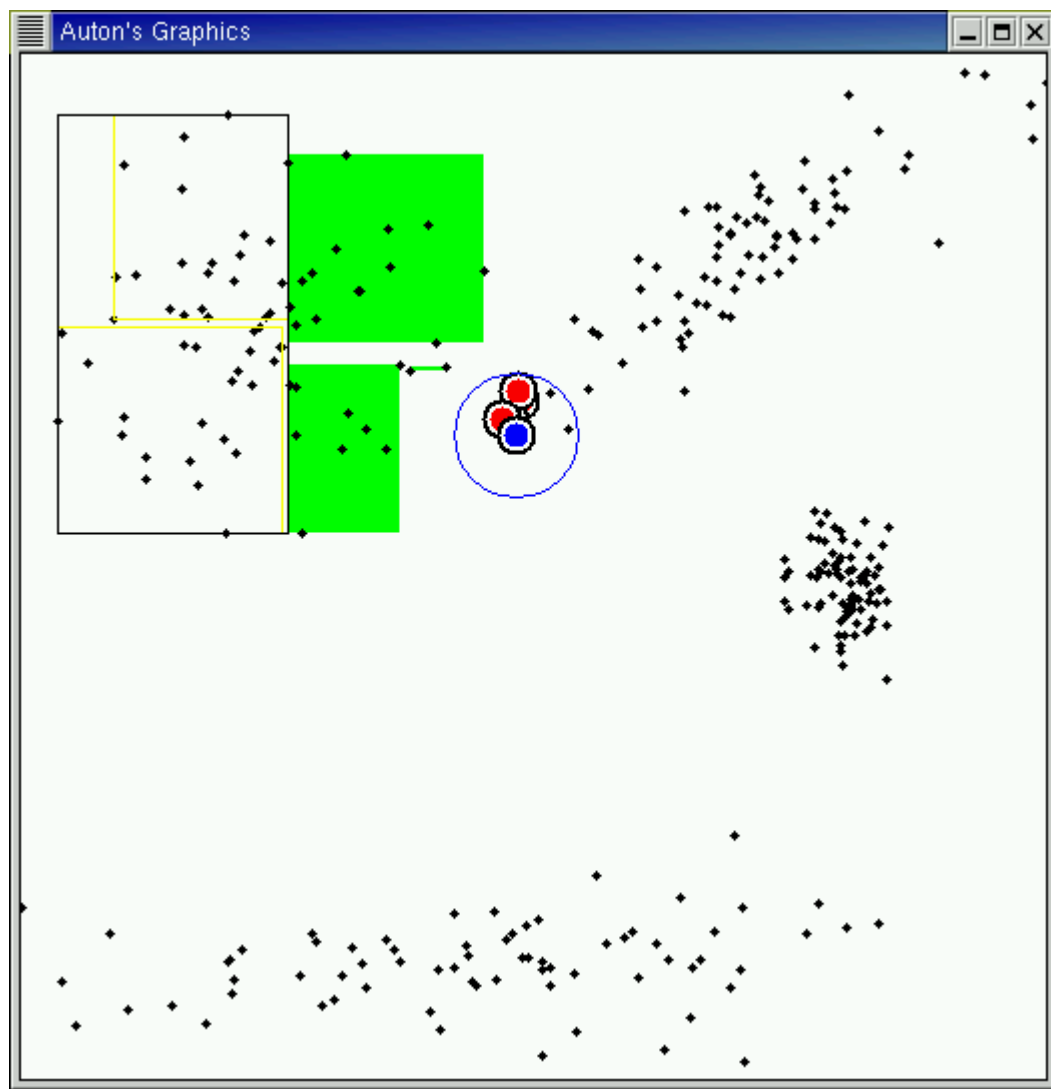


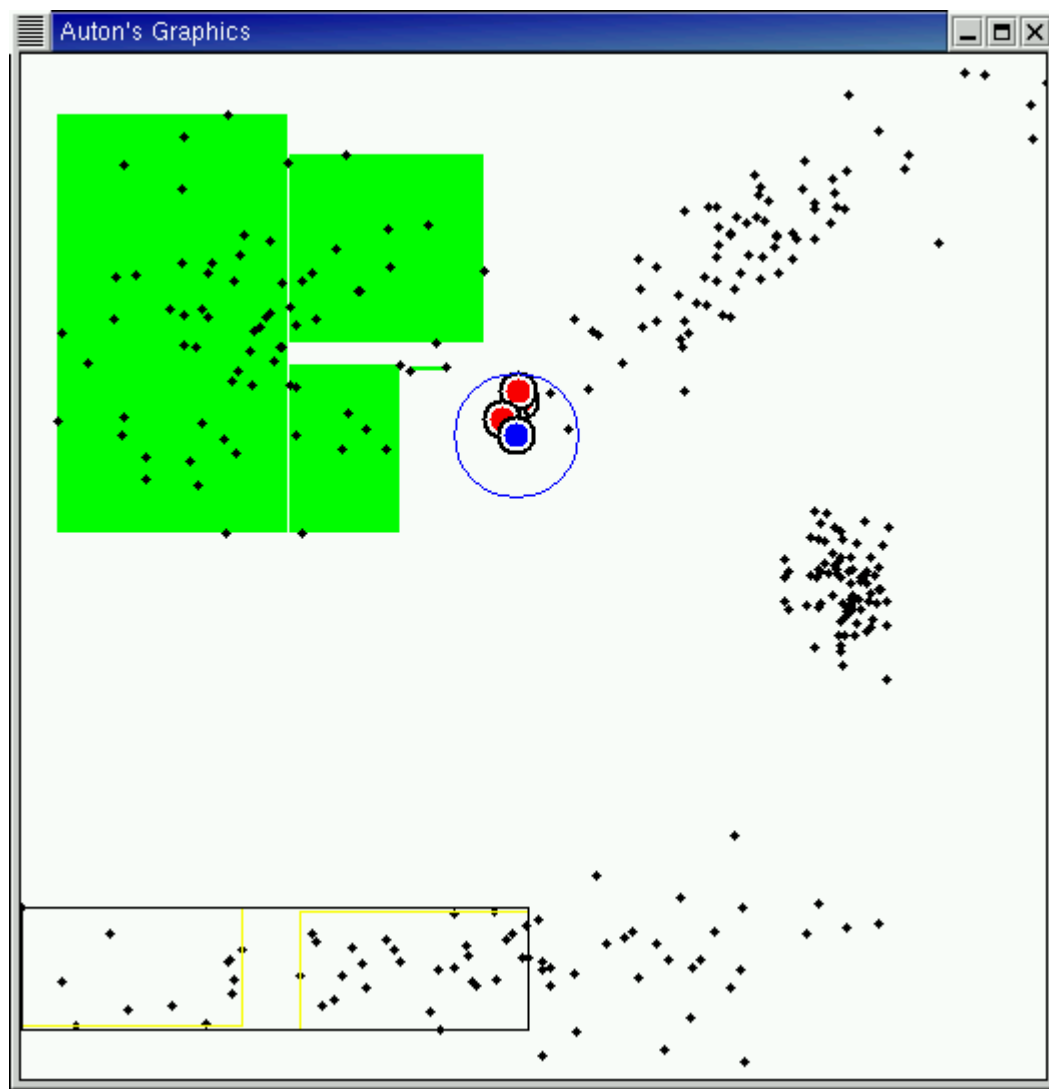


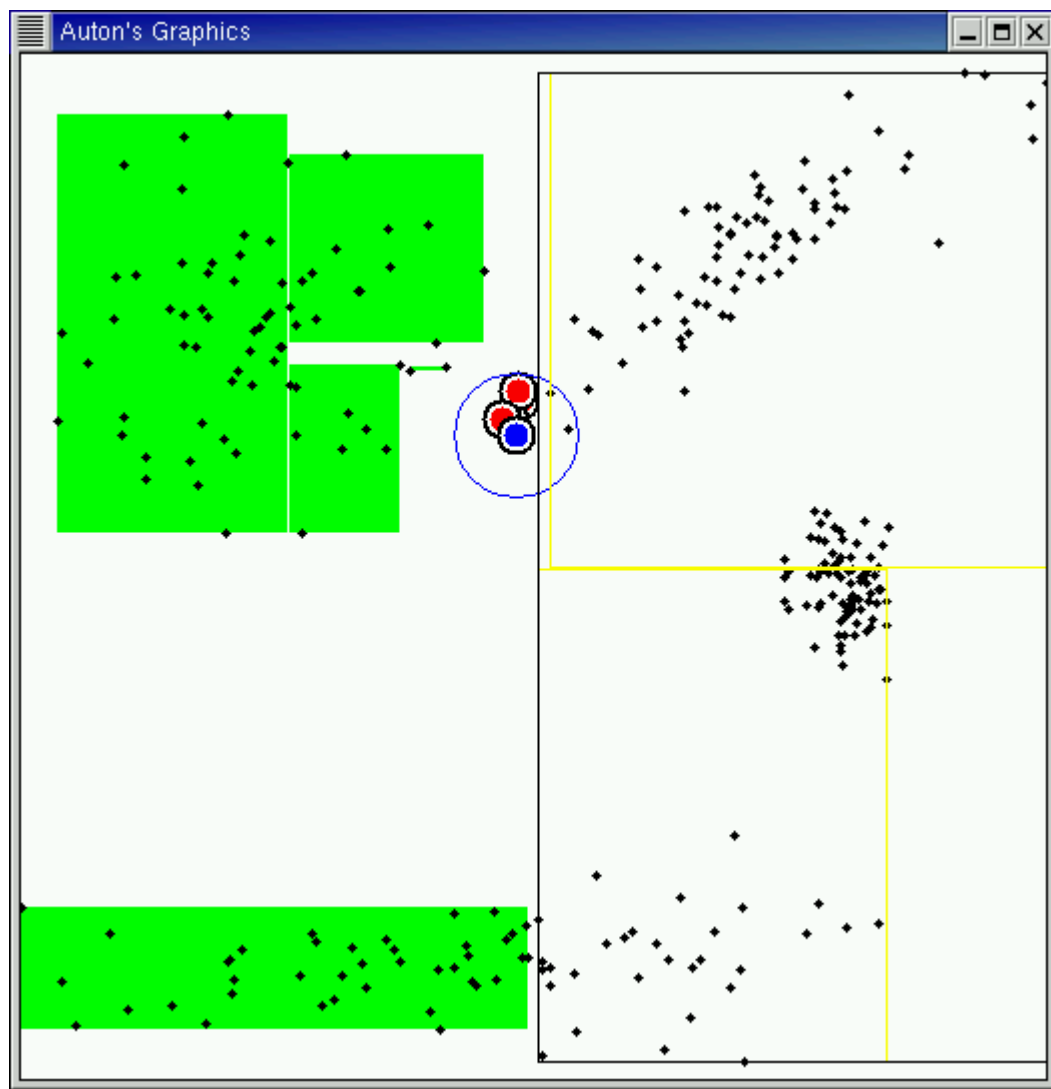


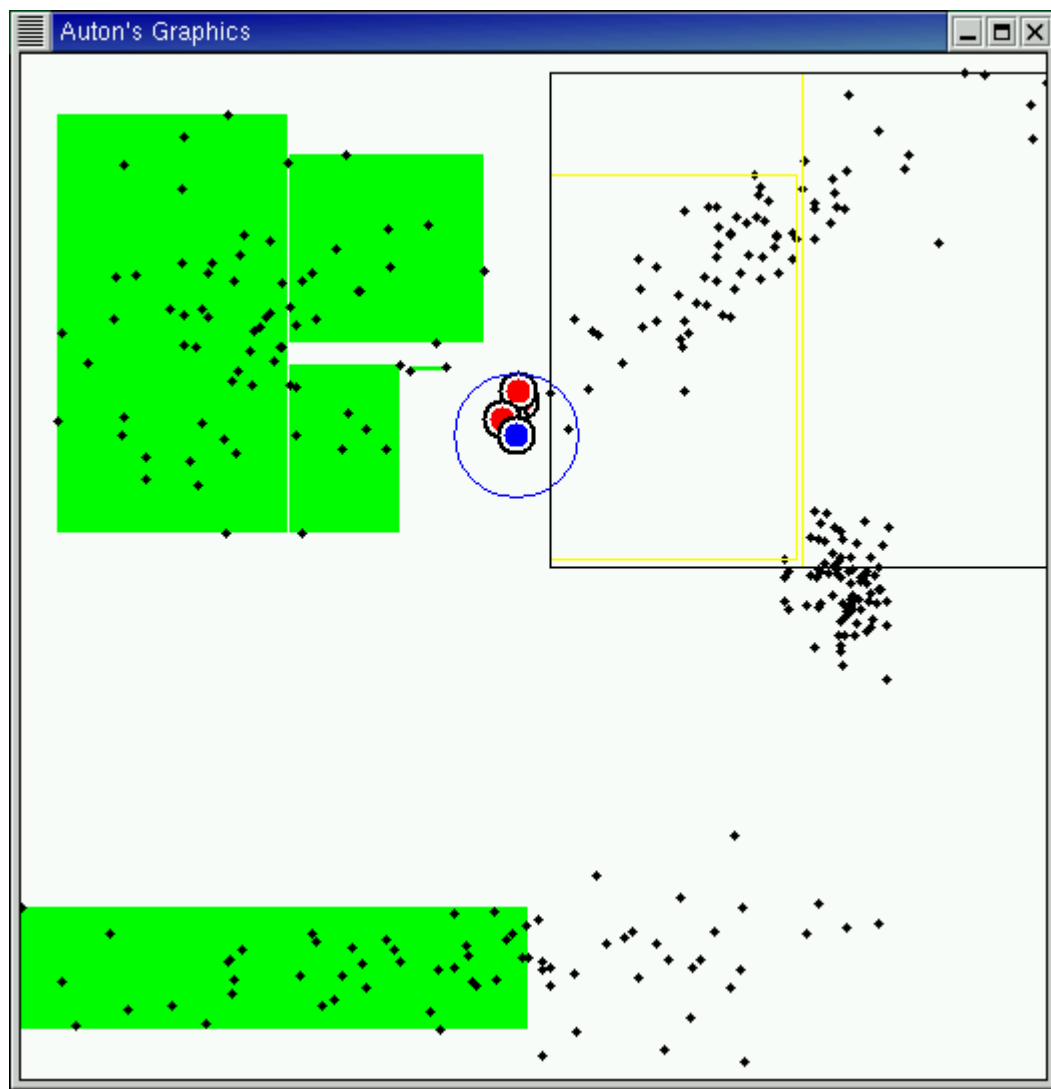


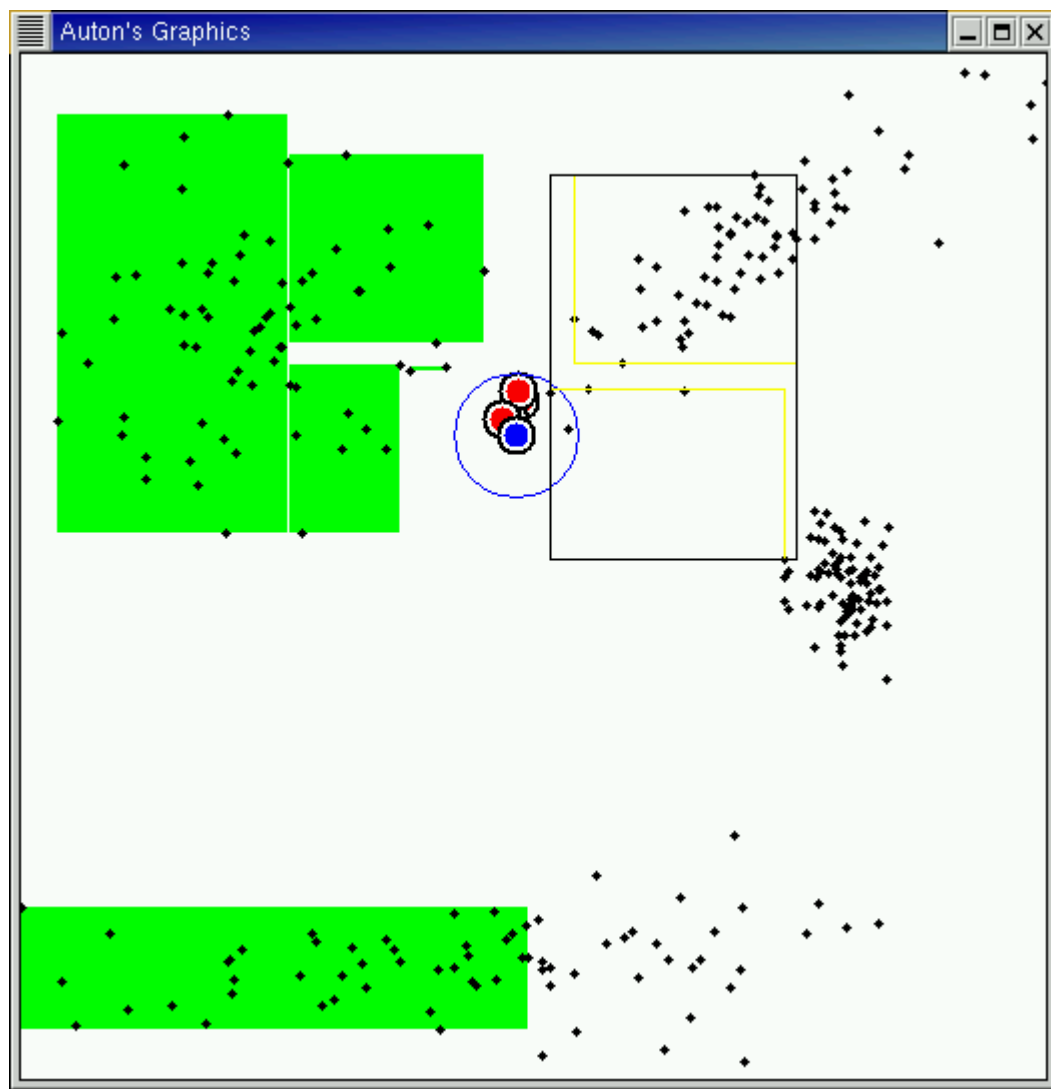


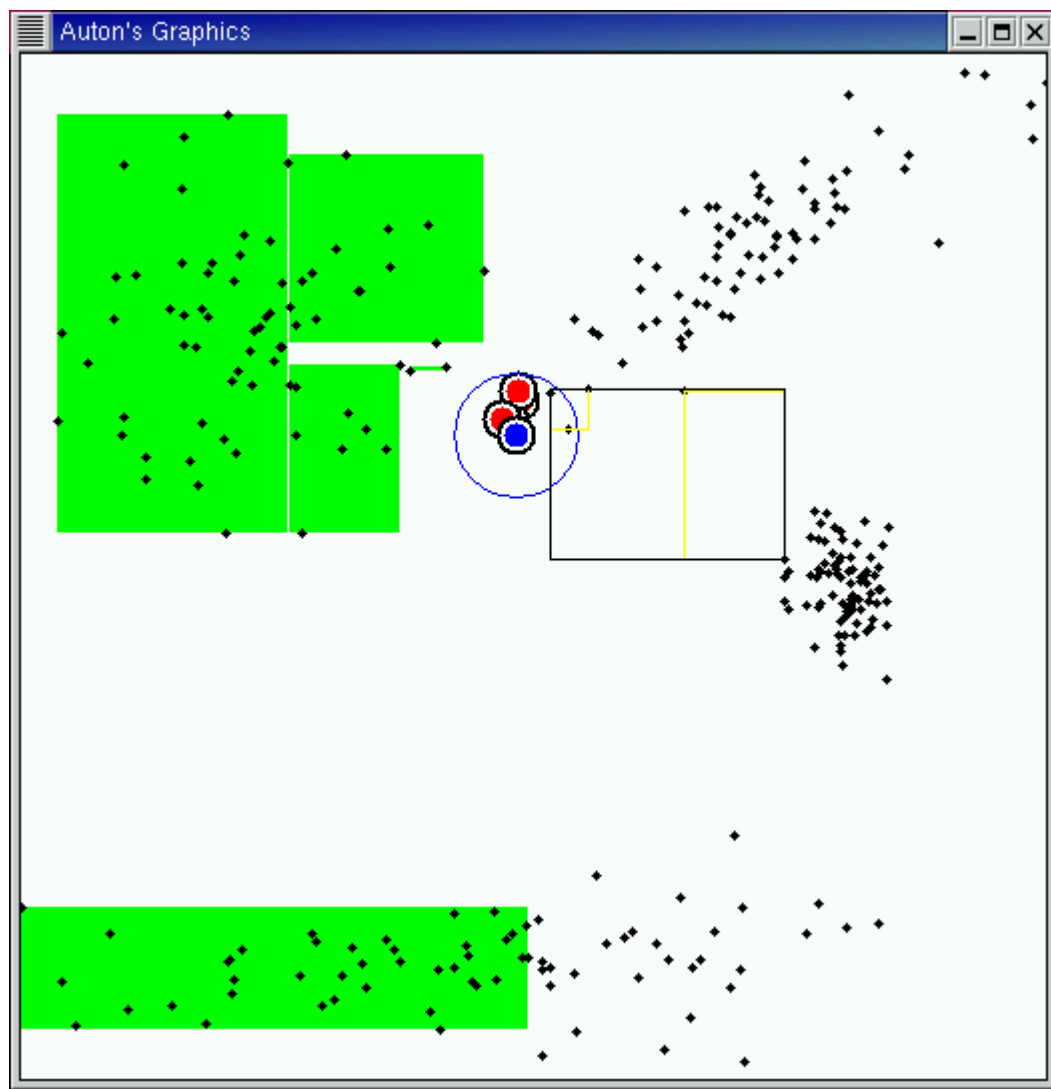


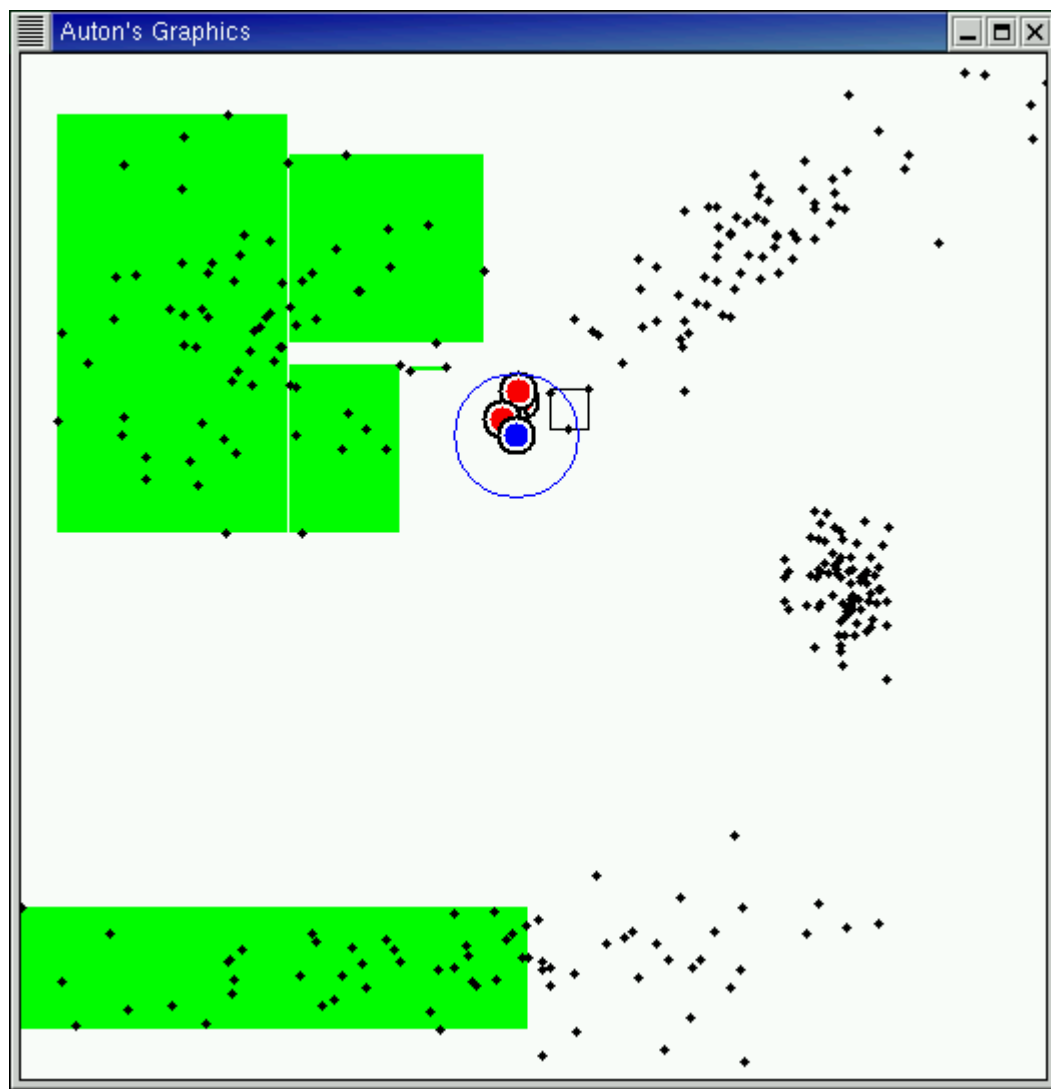


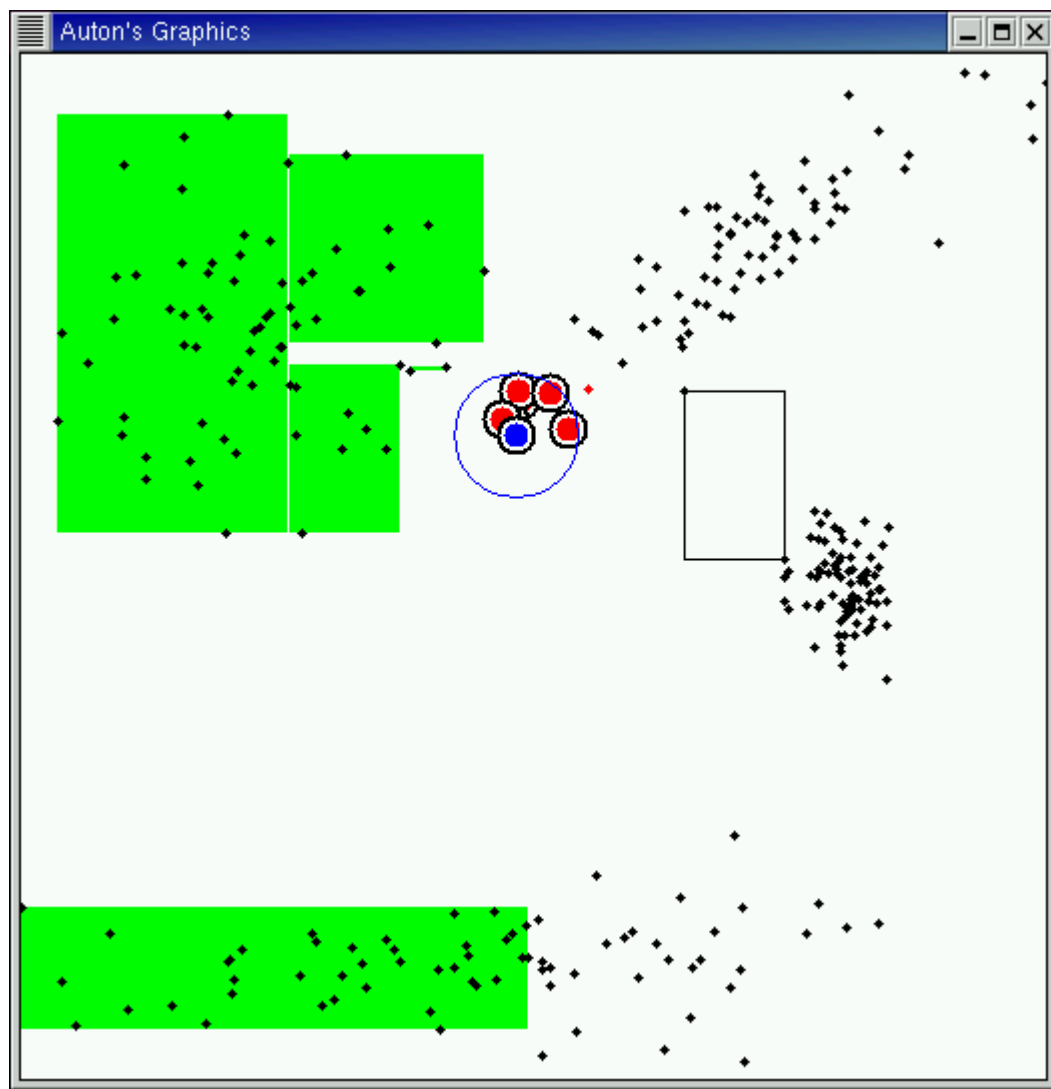


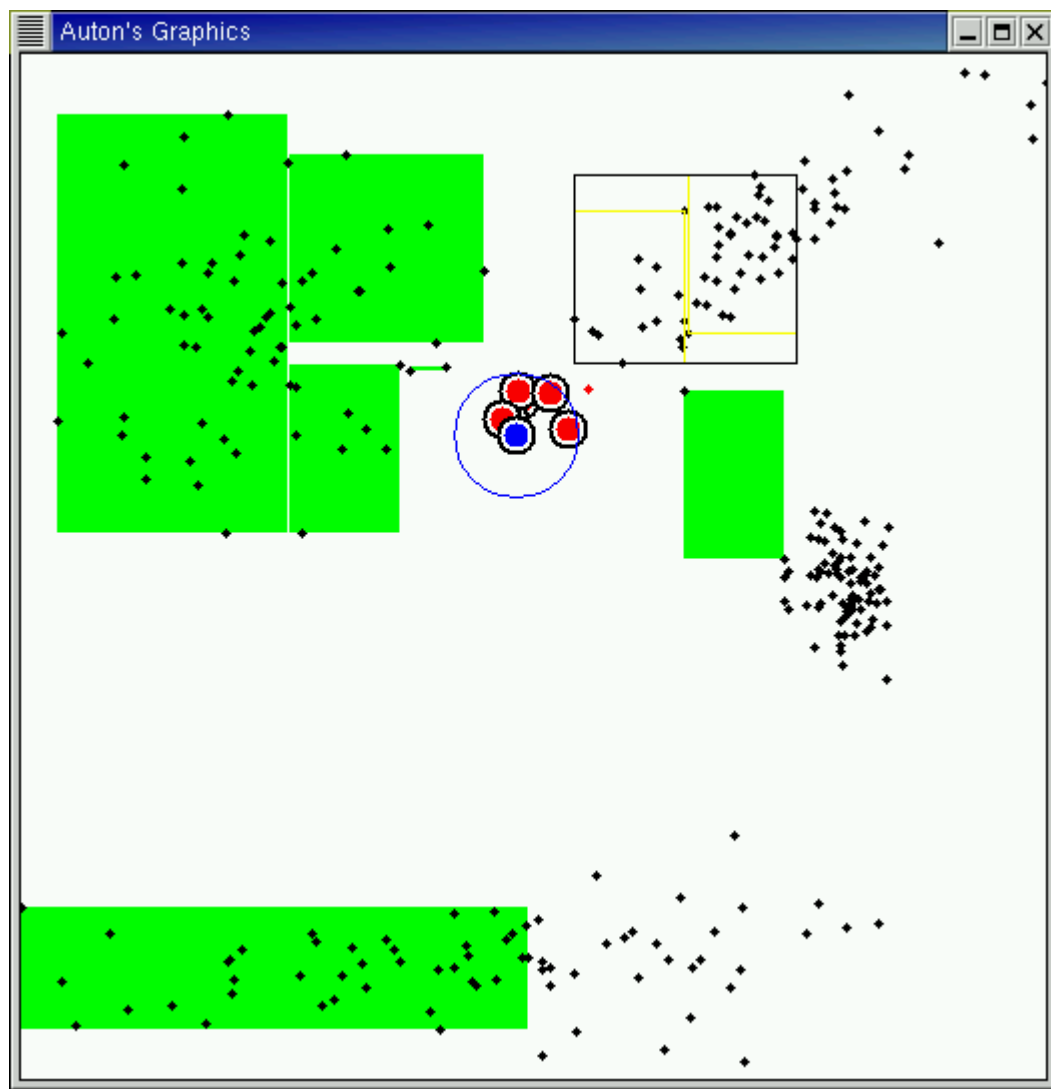


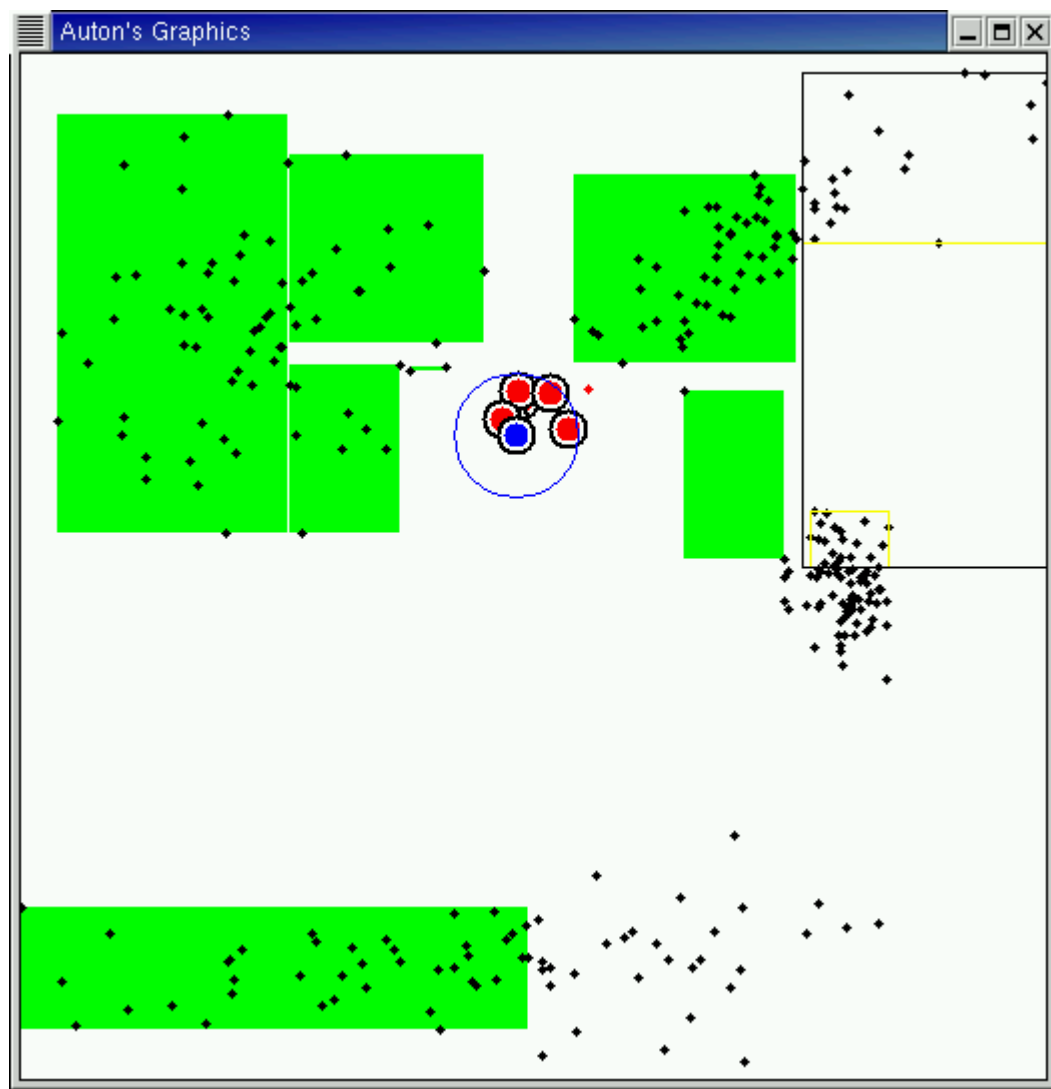


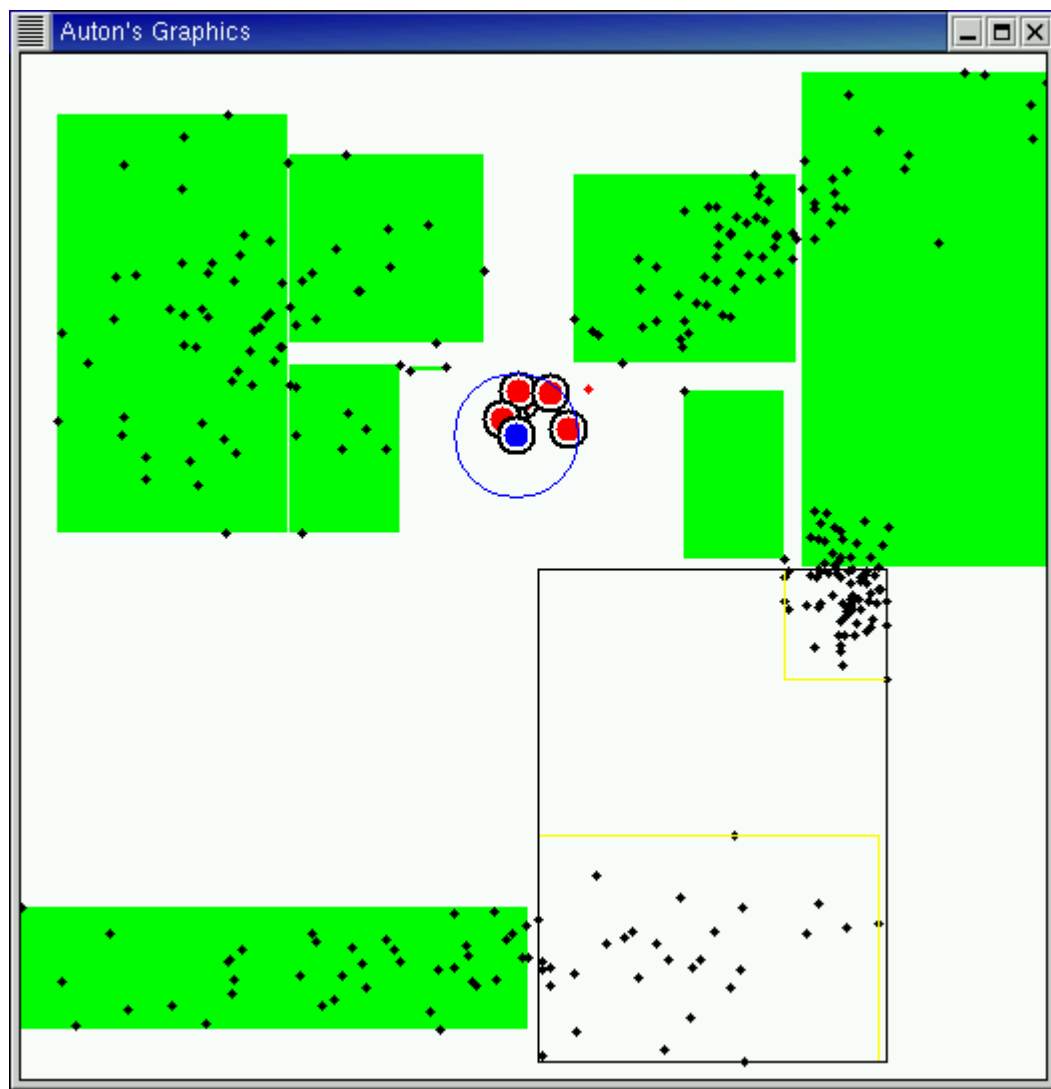


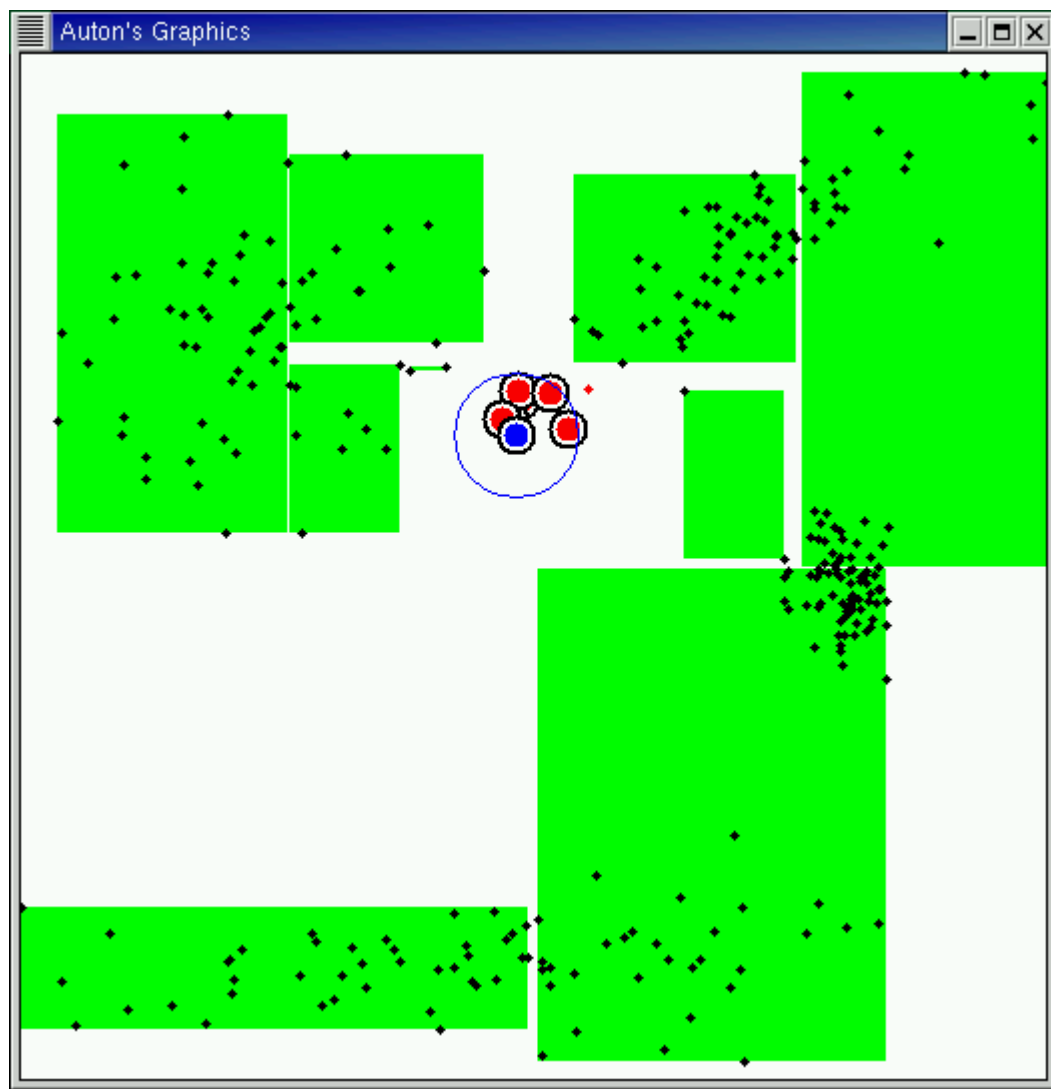




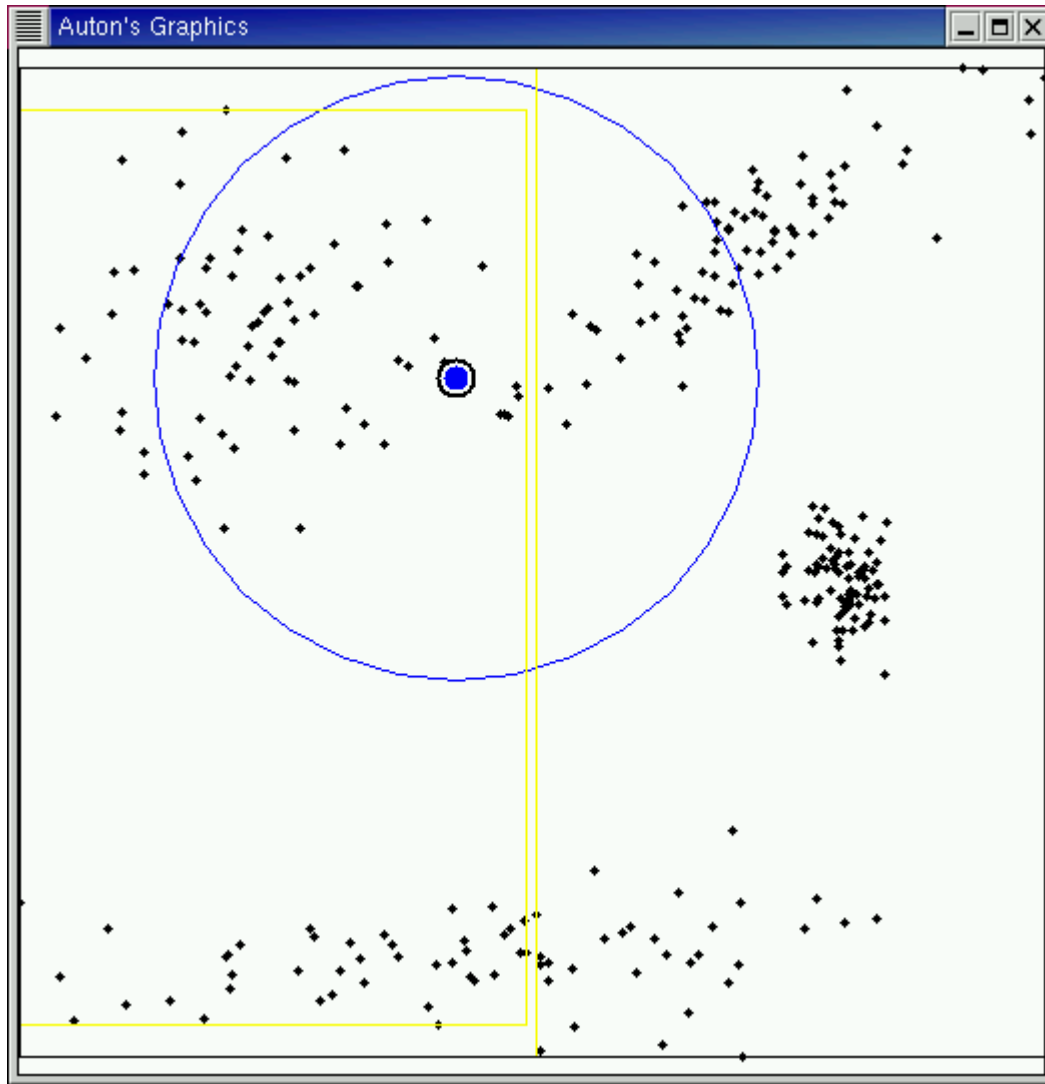




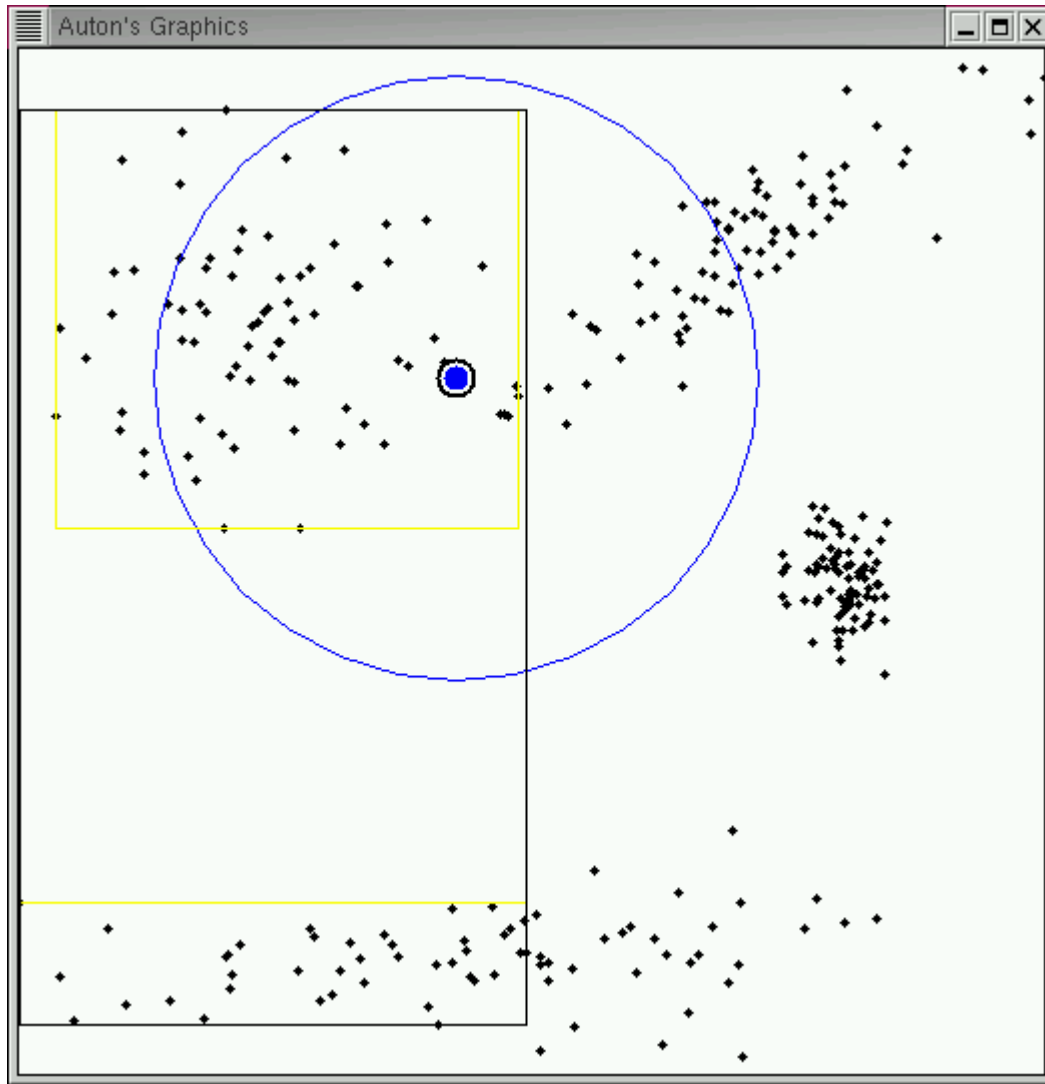




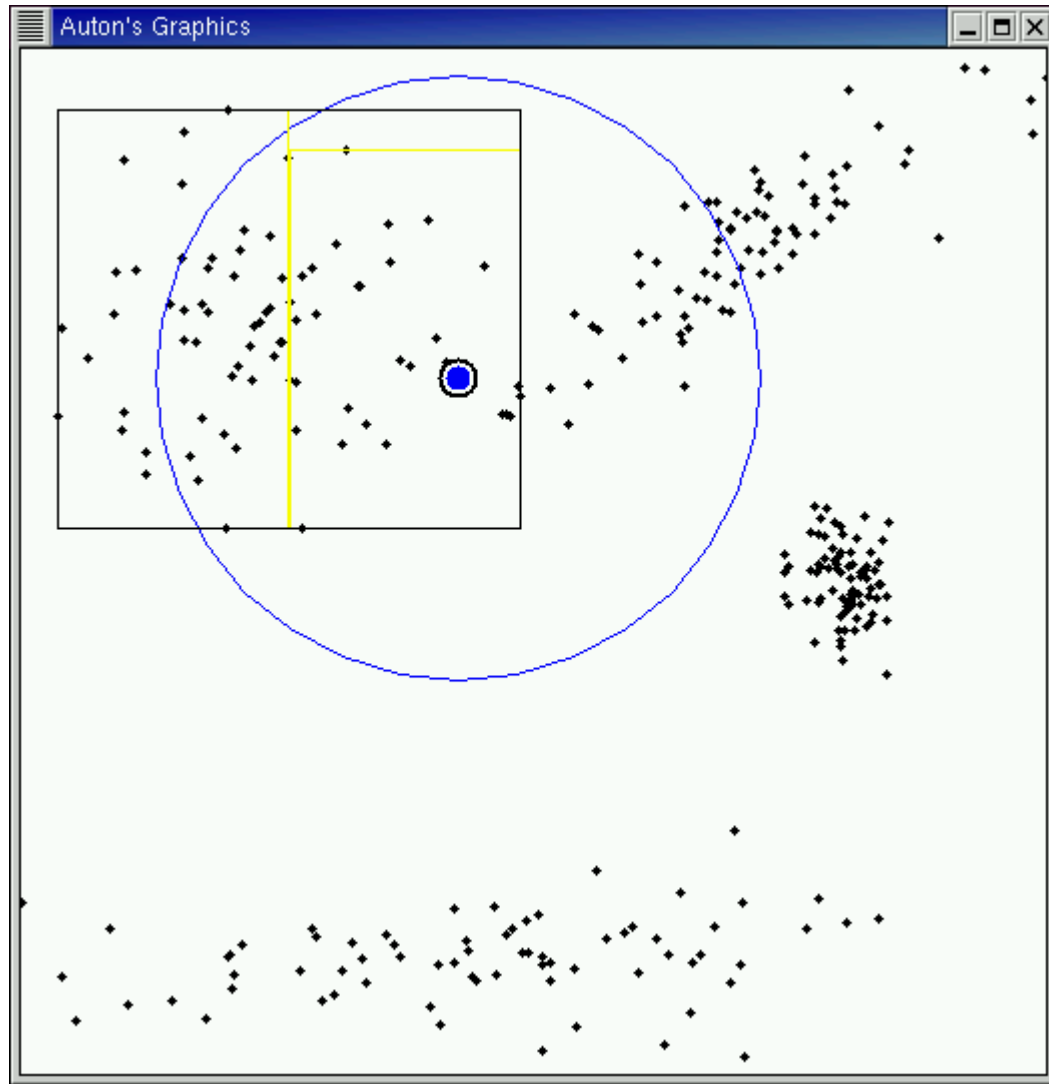
Range Count



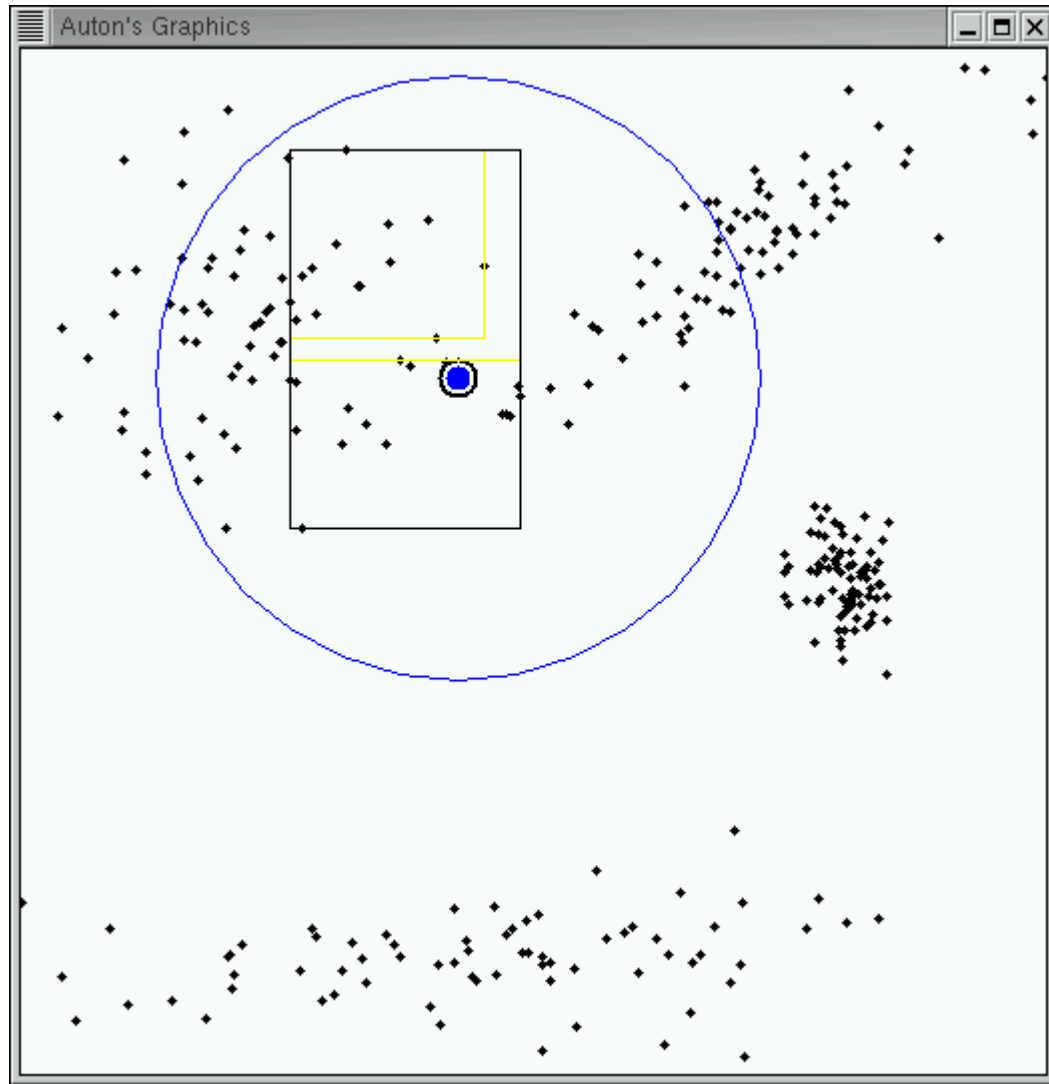
Range Count



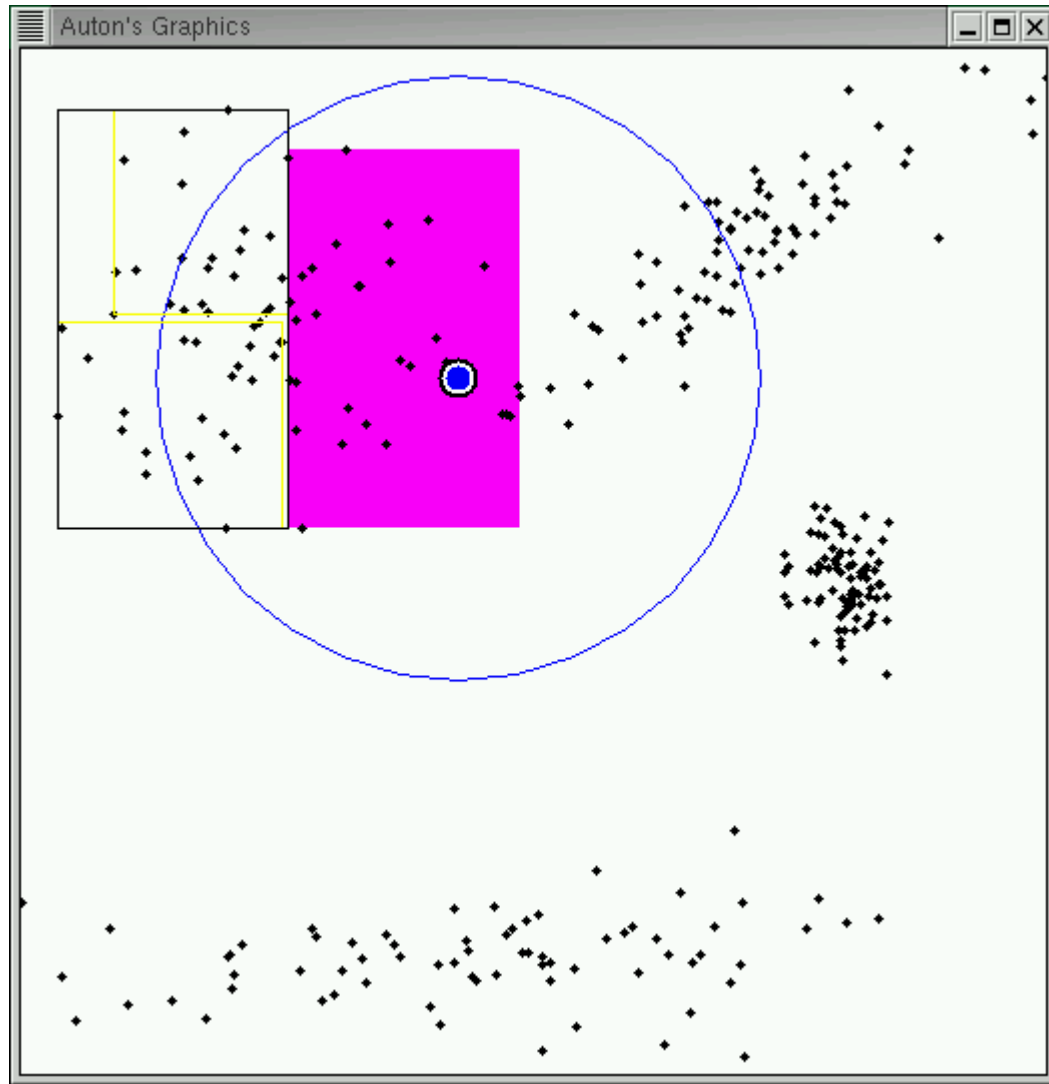
Range Count



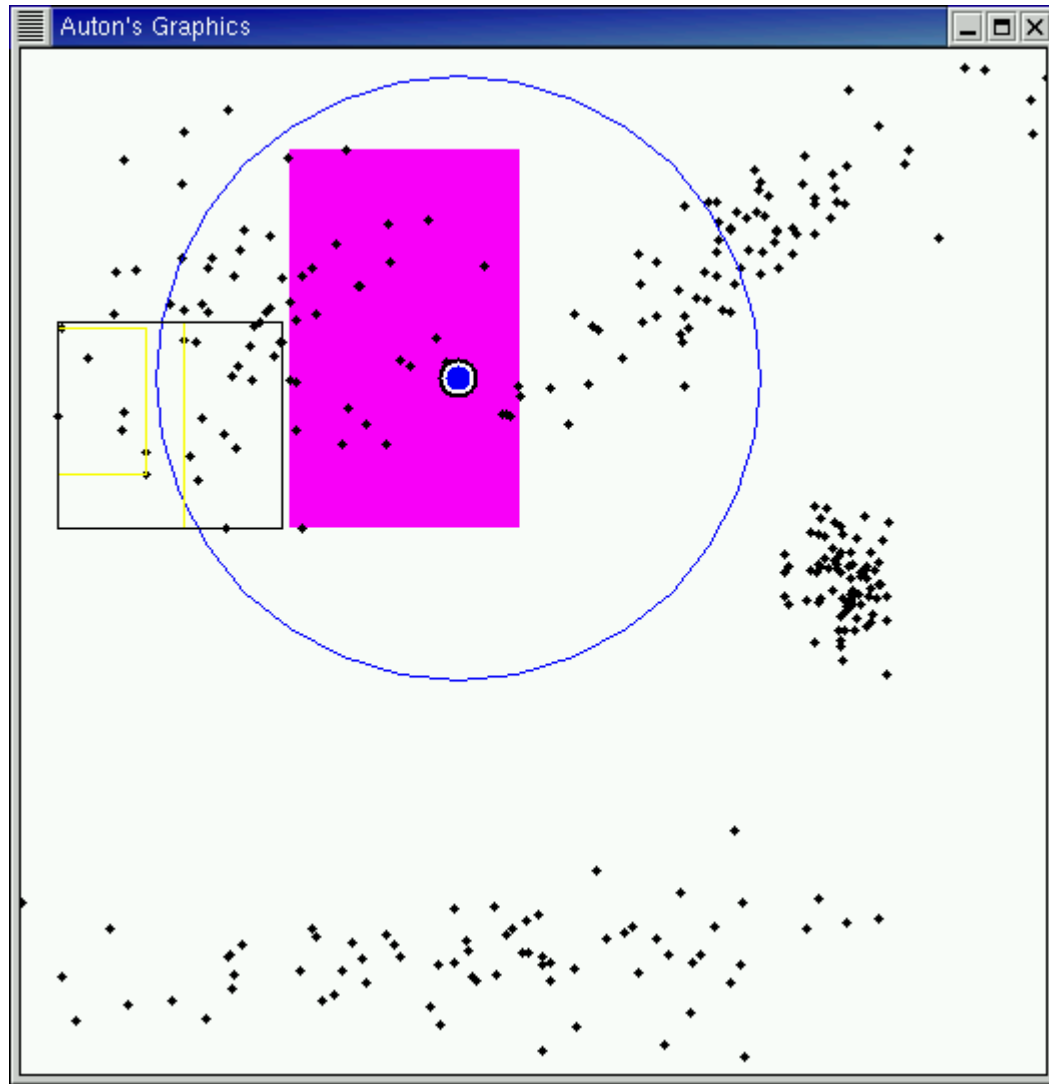
Range Count



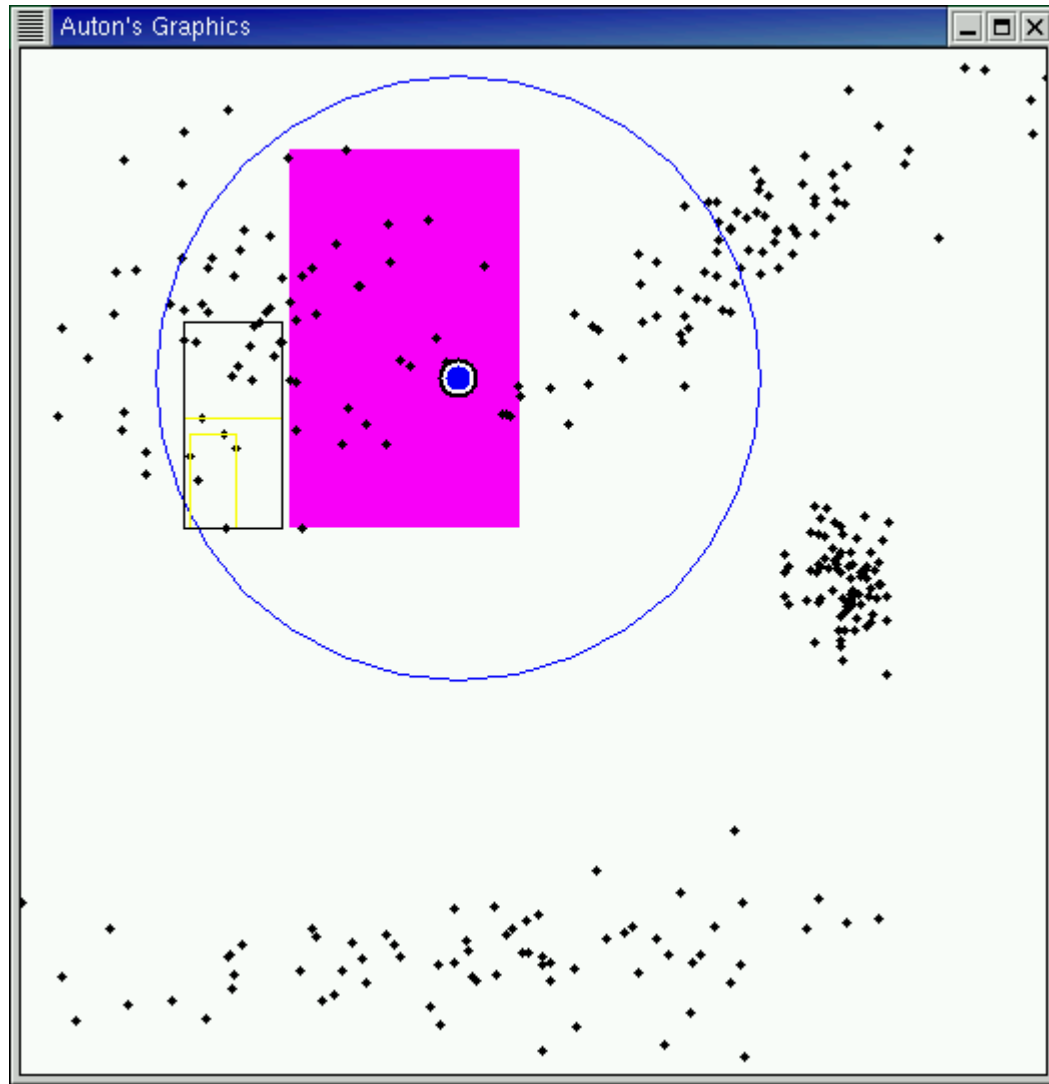
Range Count



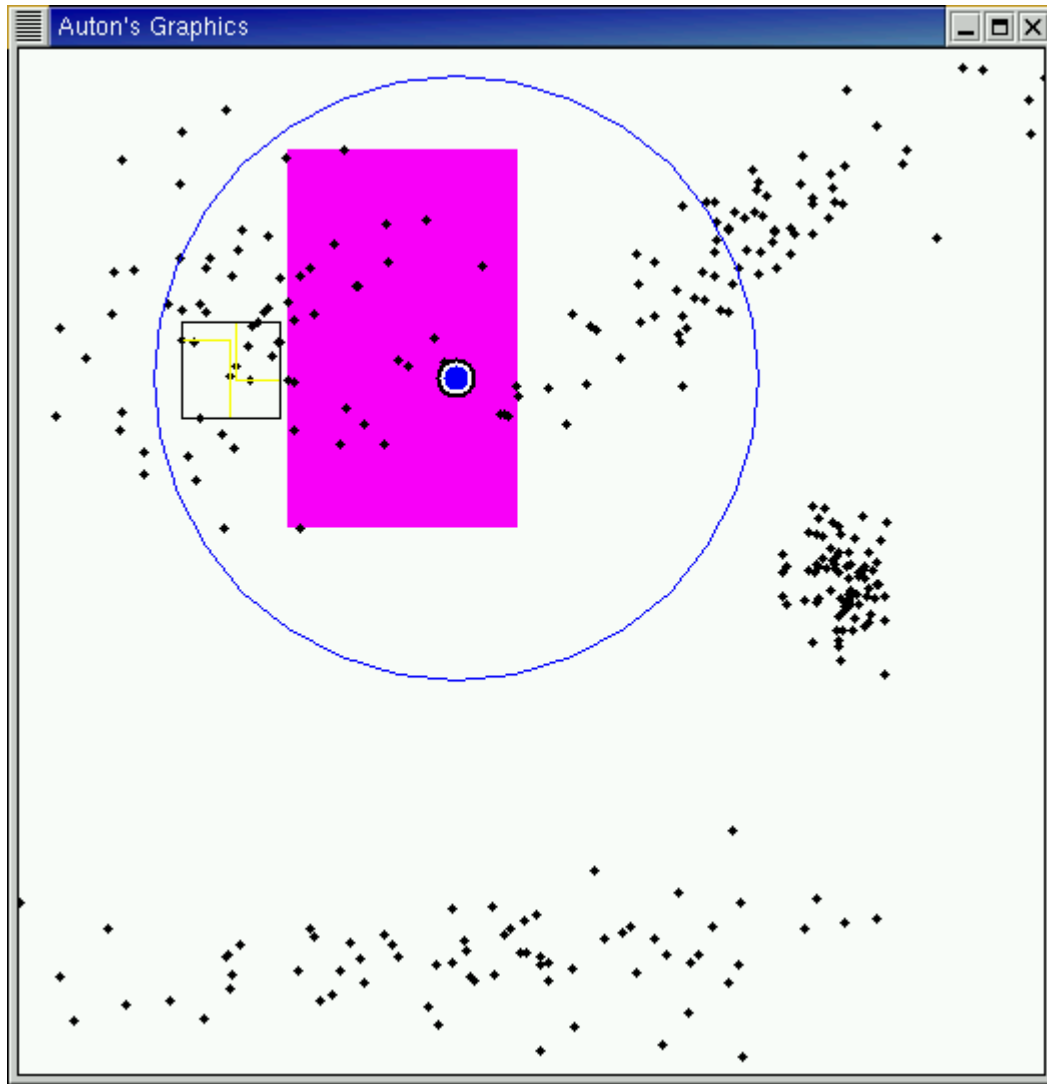
Range Count



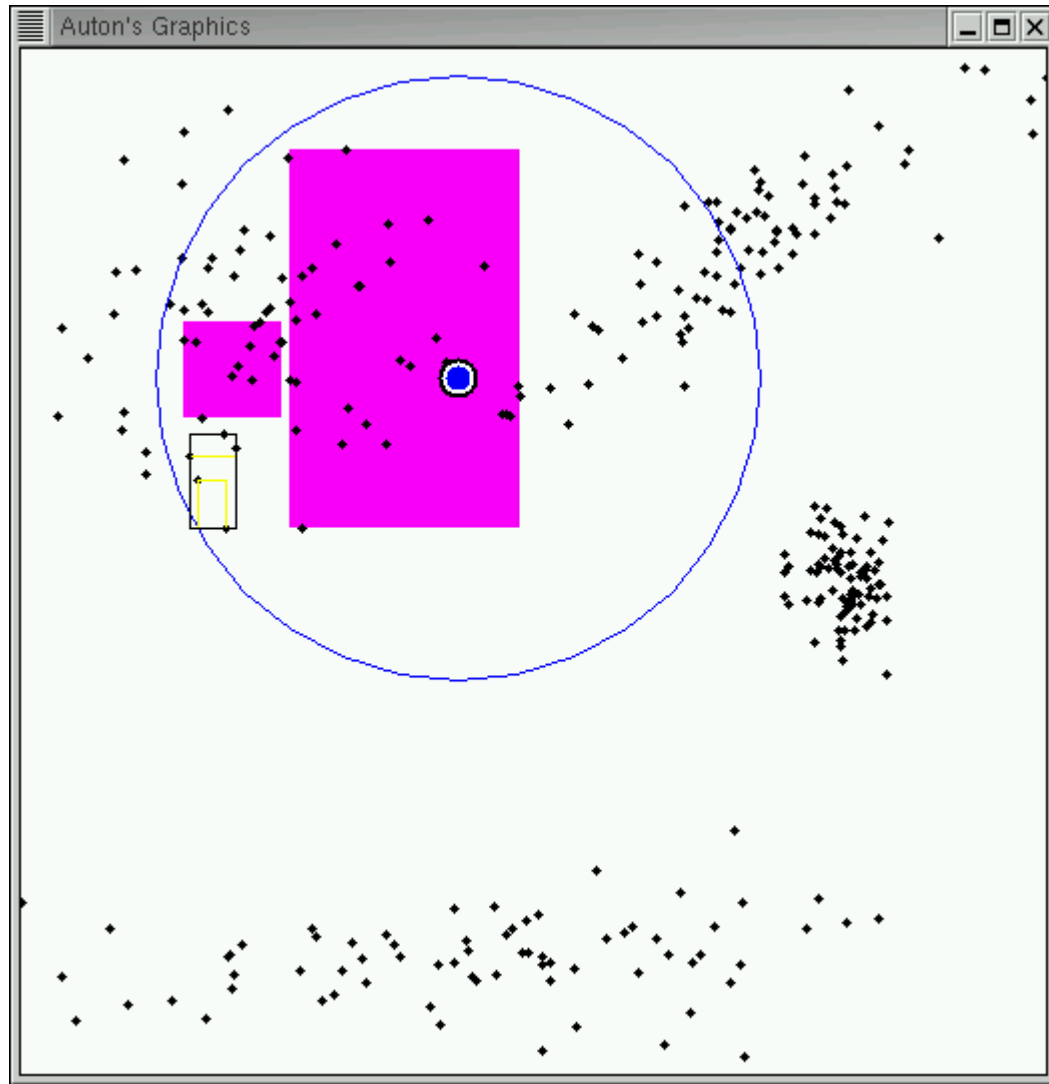
Range Count



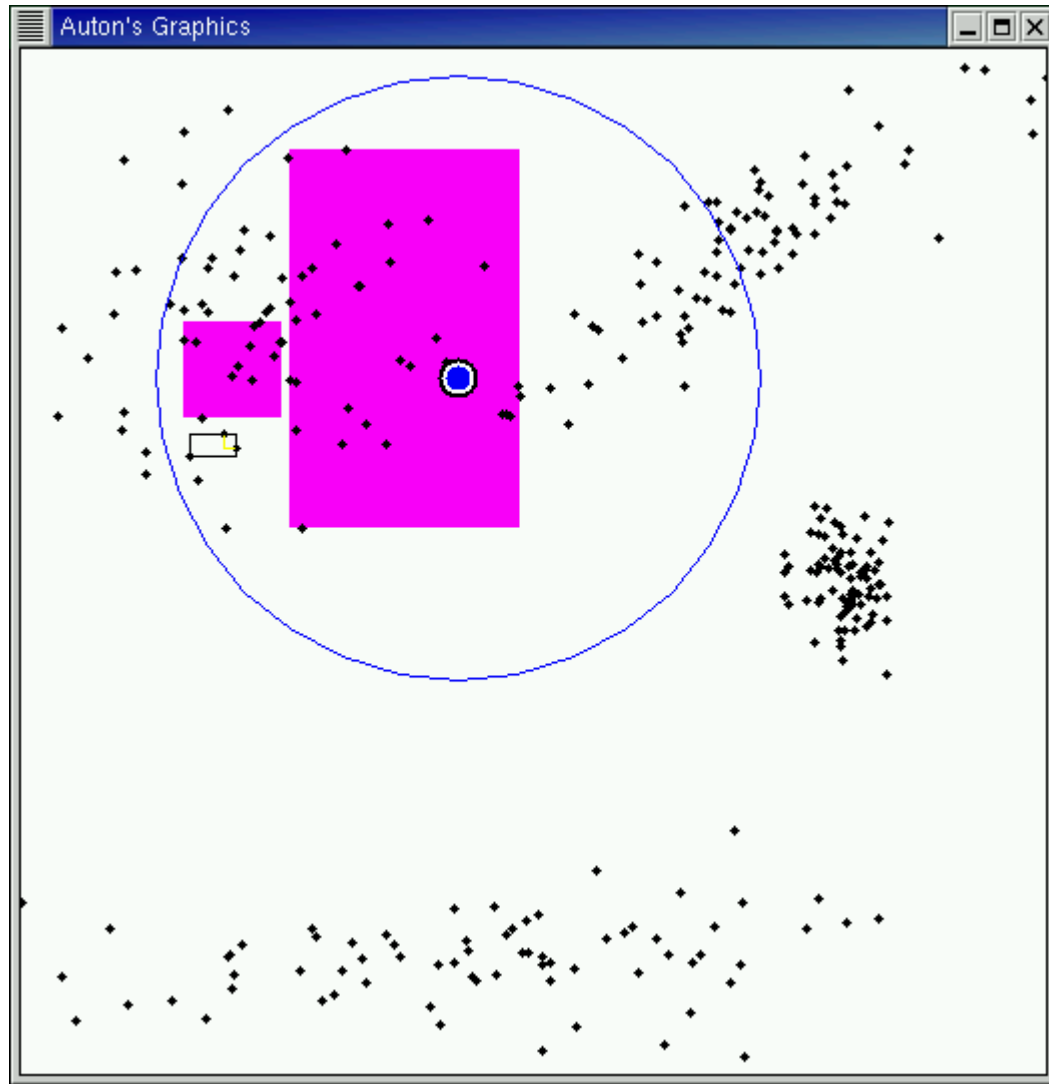
Range Count



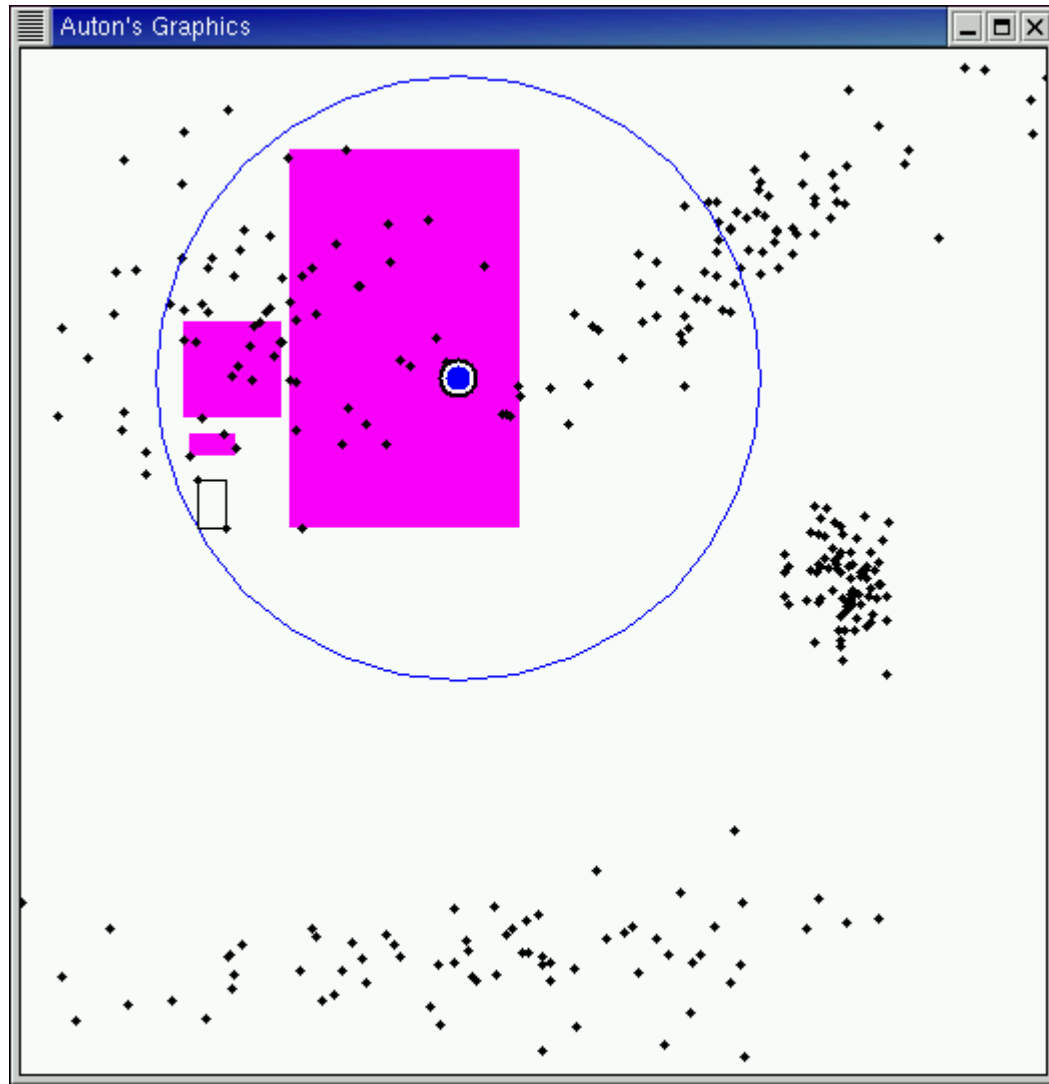
Range Count



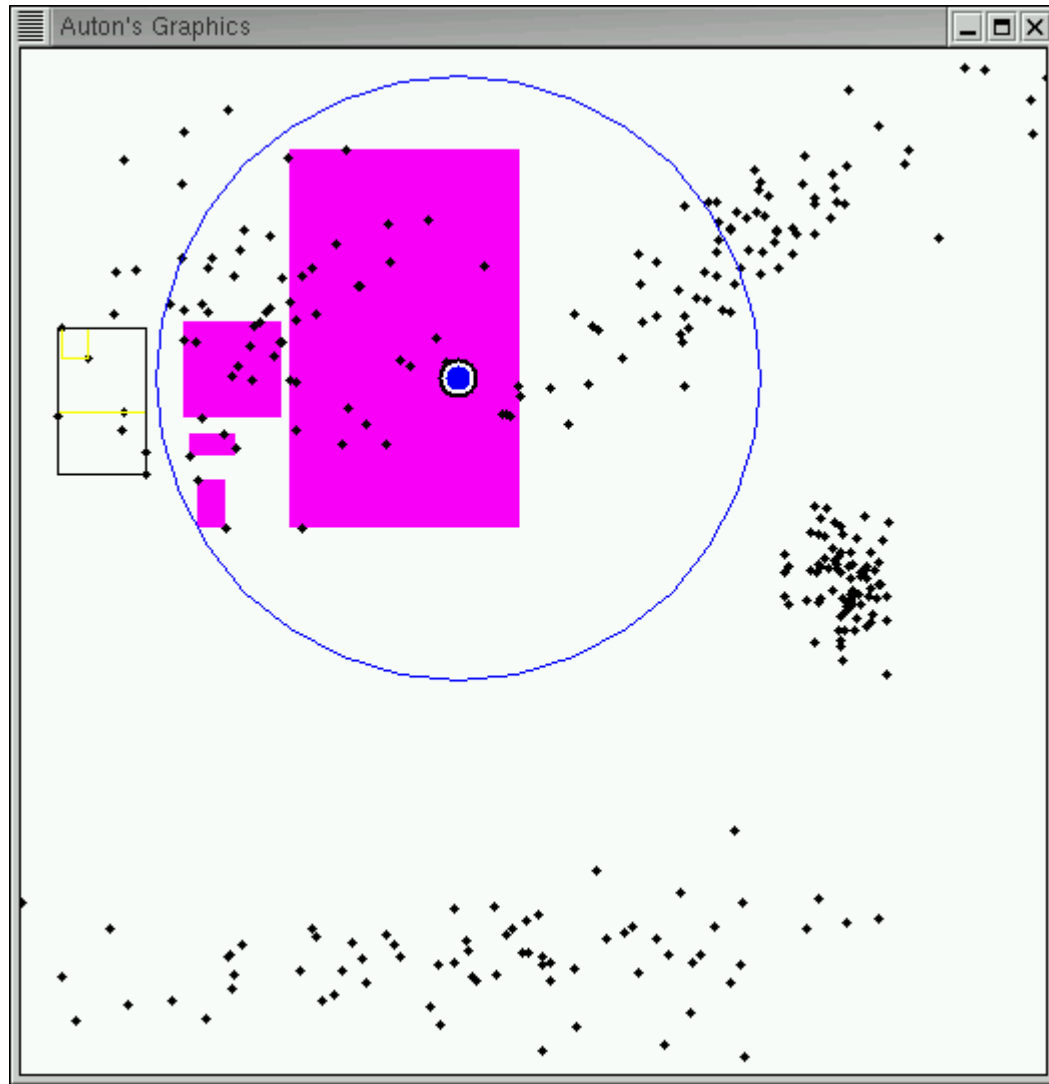
Range Count



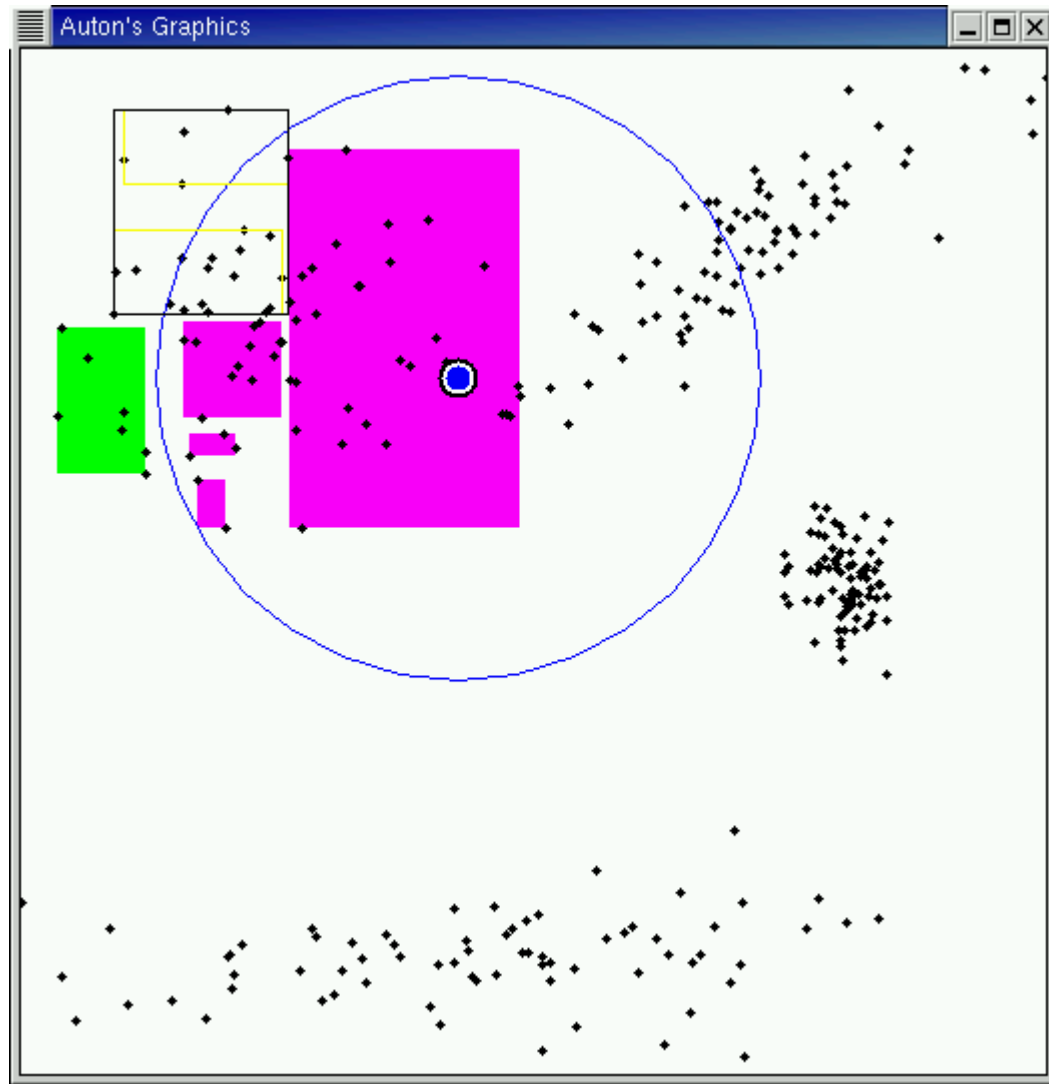
Range Count



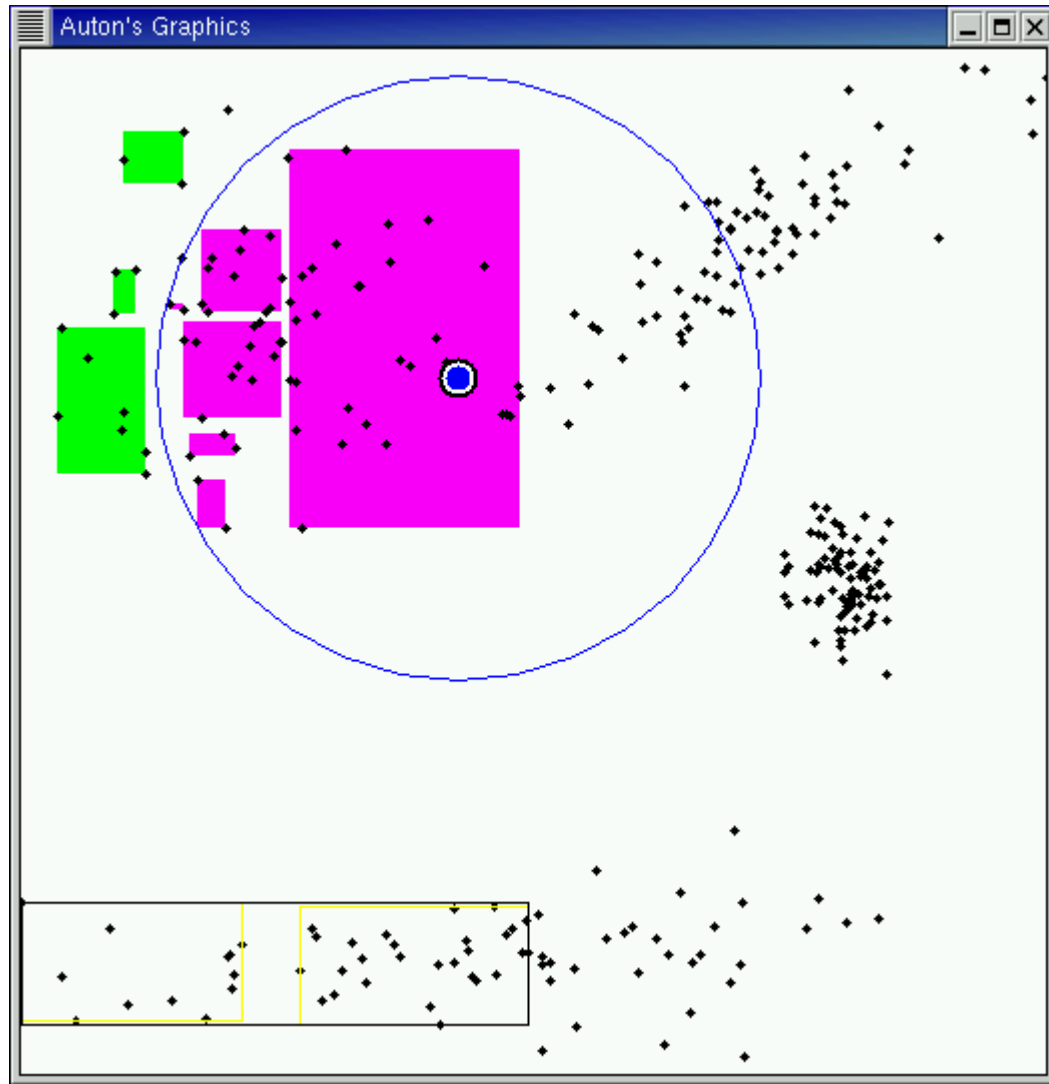
Range Count



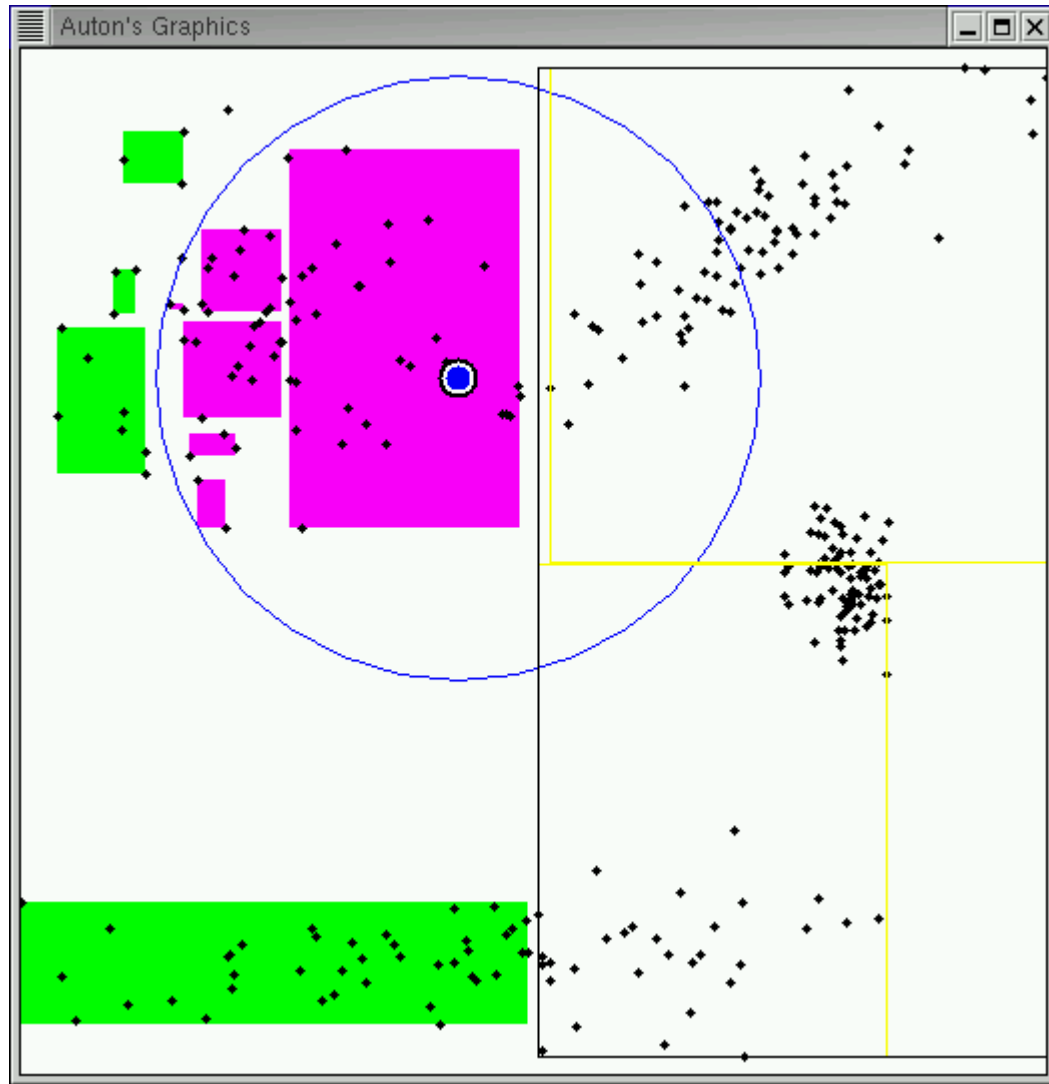
Range Count



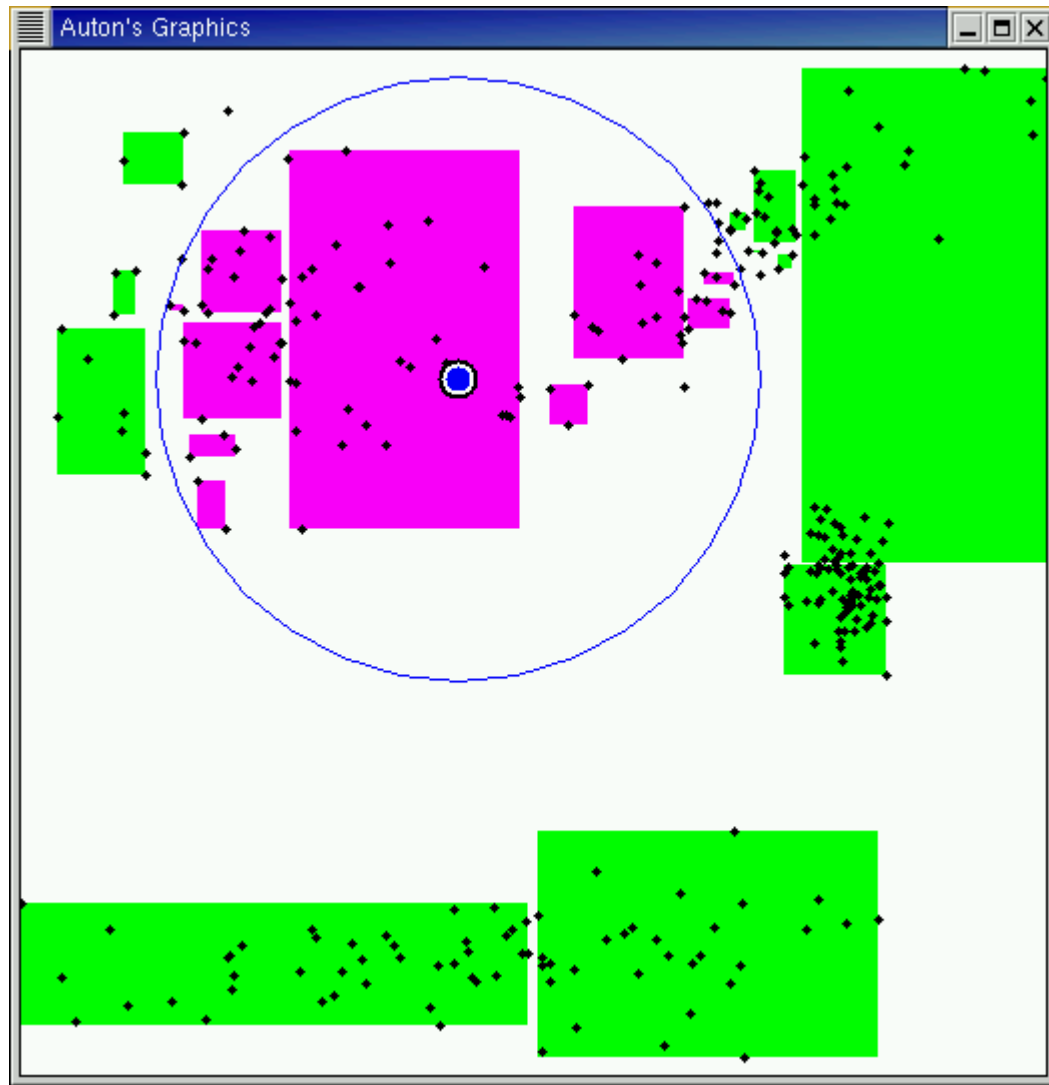
Range Count



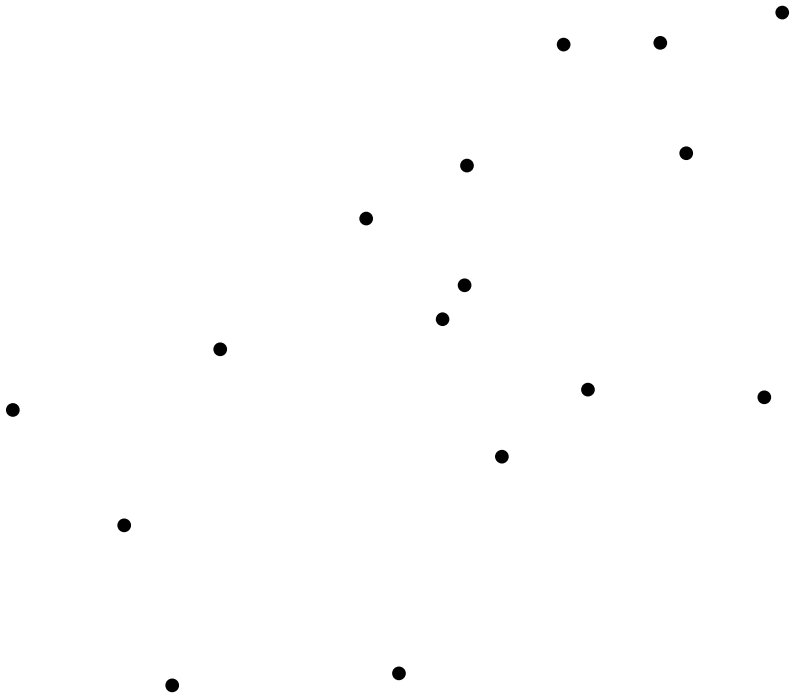
Range Count

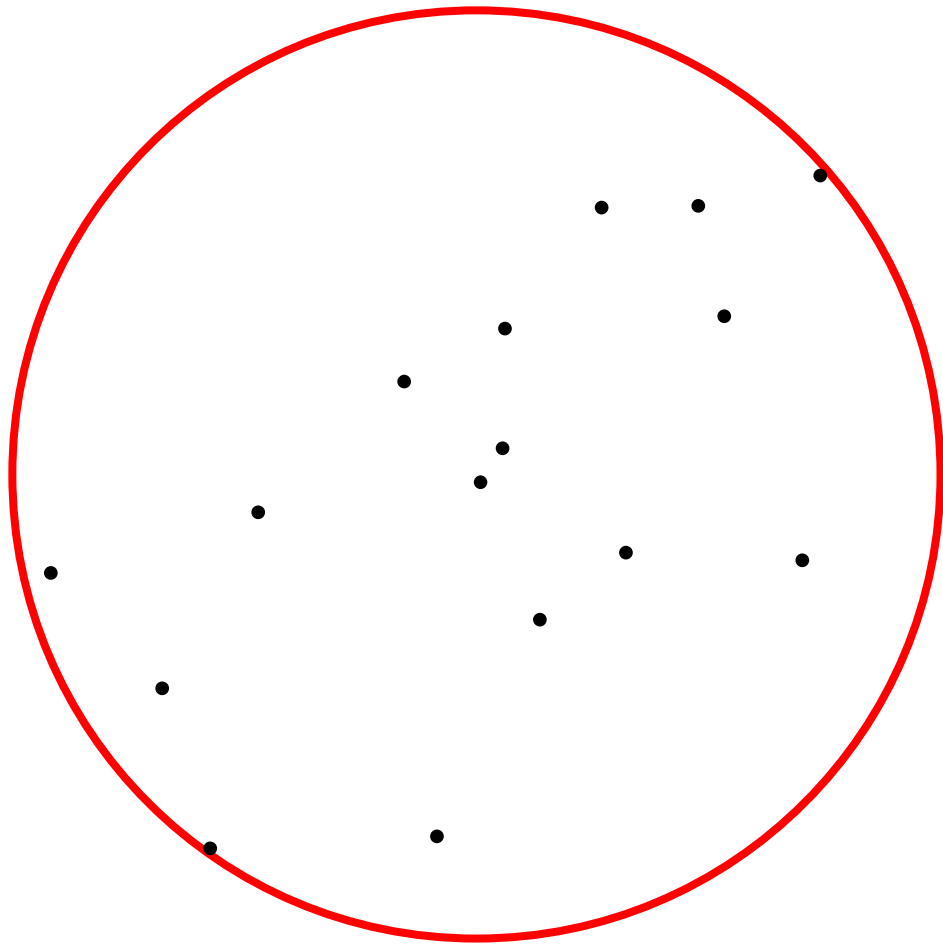


Range Count



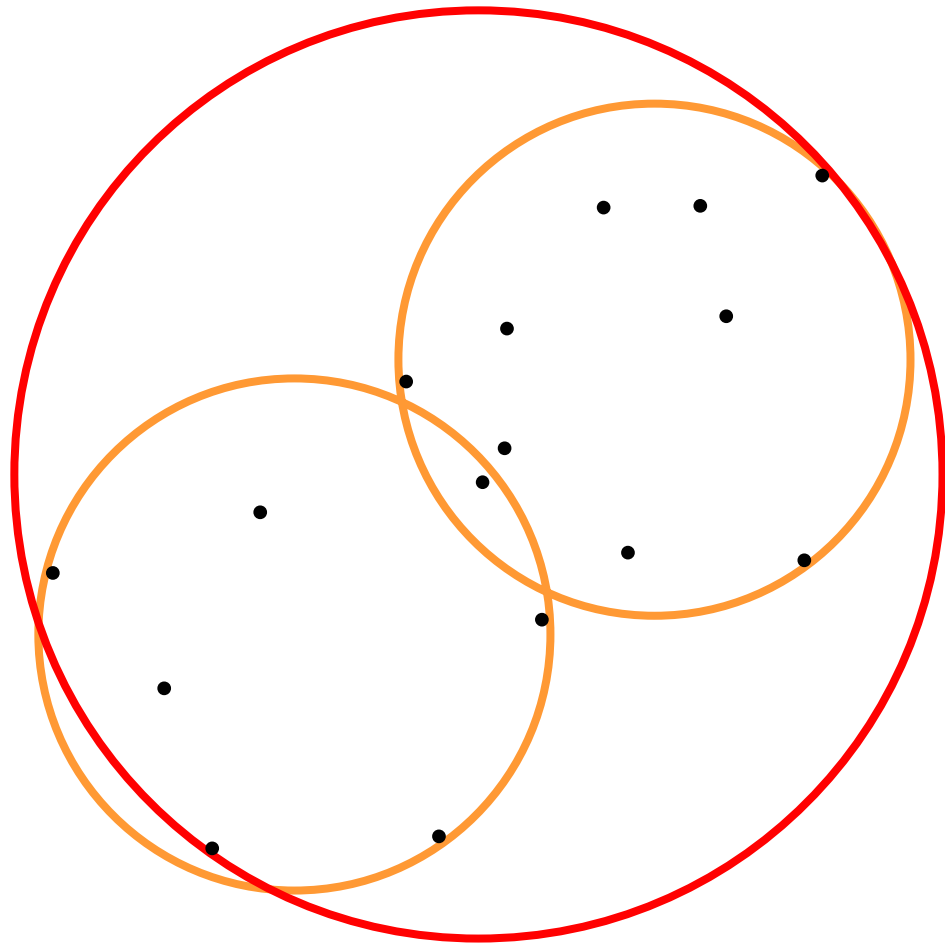
A Set of Points in a metric space



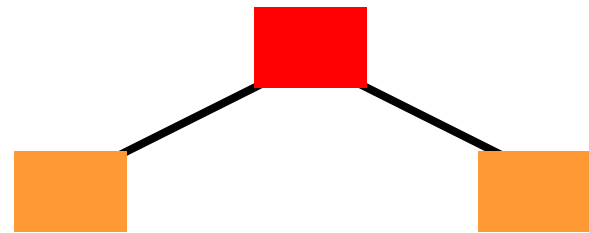


Ball Tree root
node

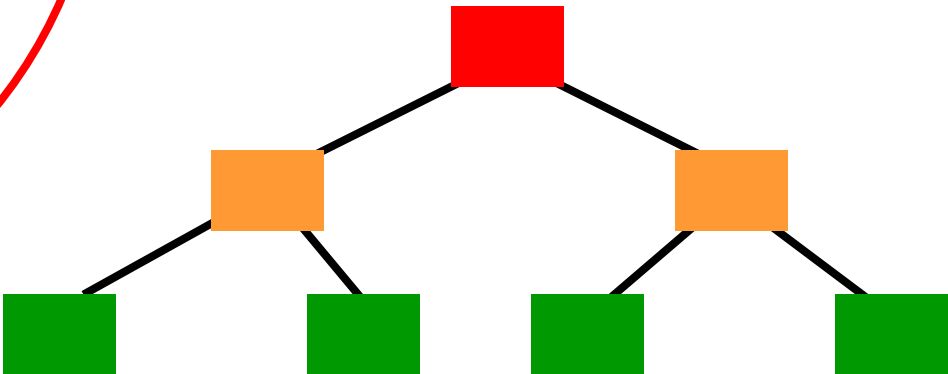
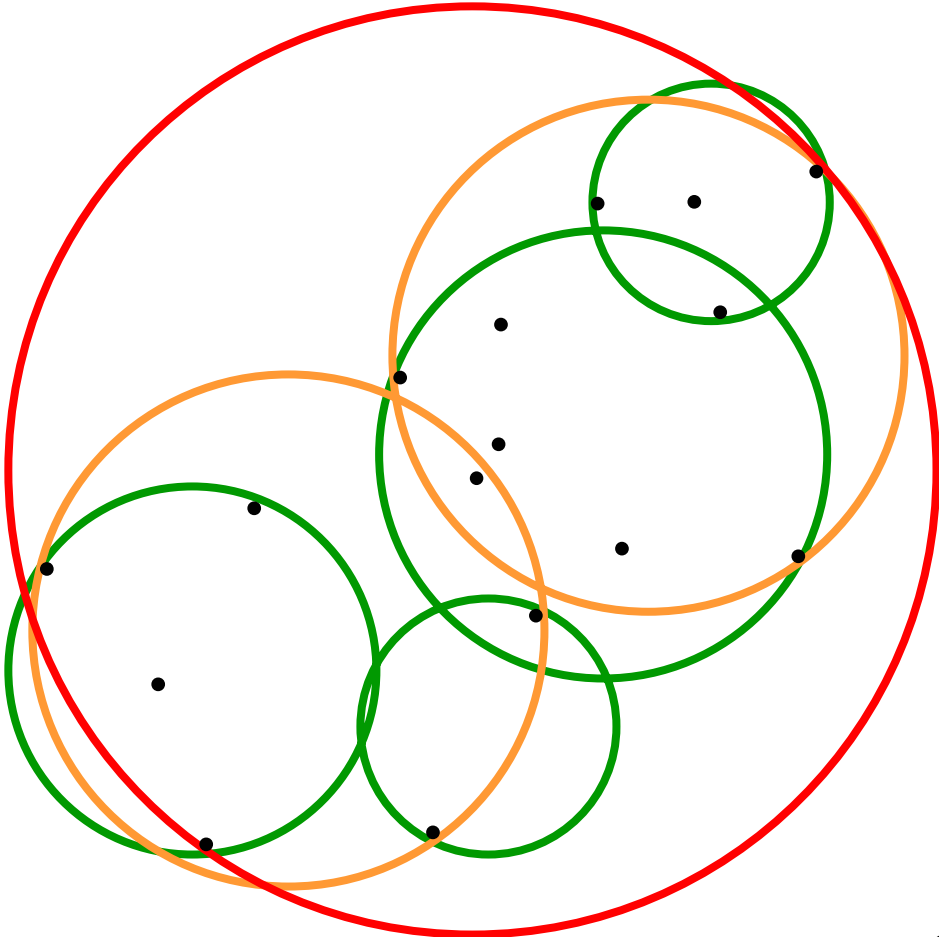




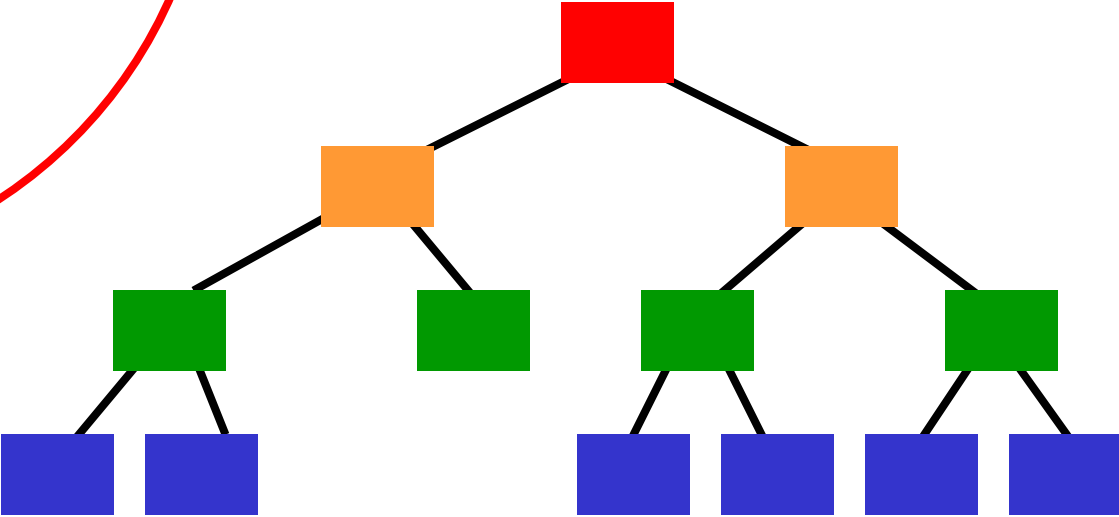
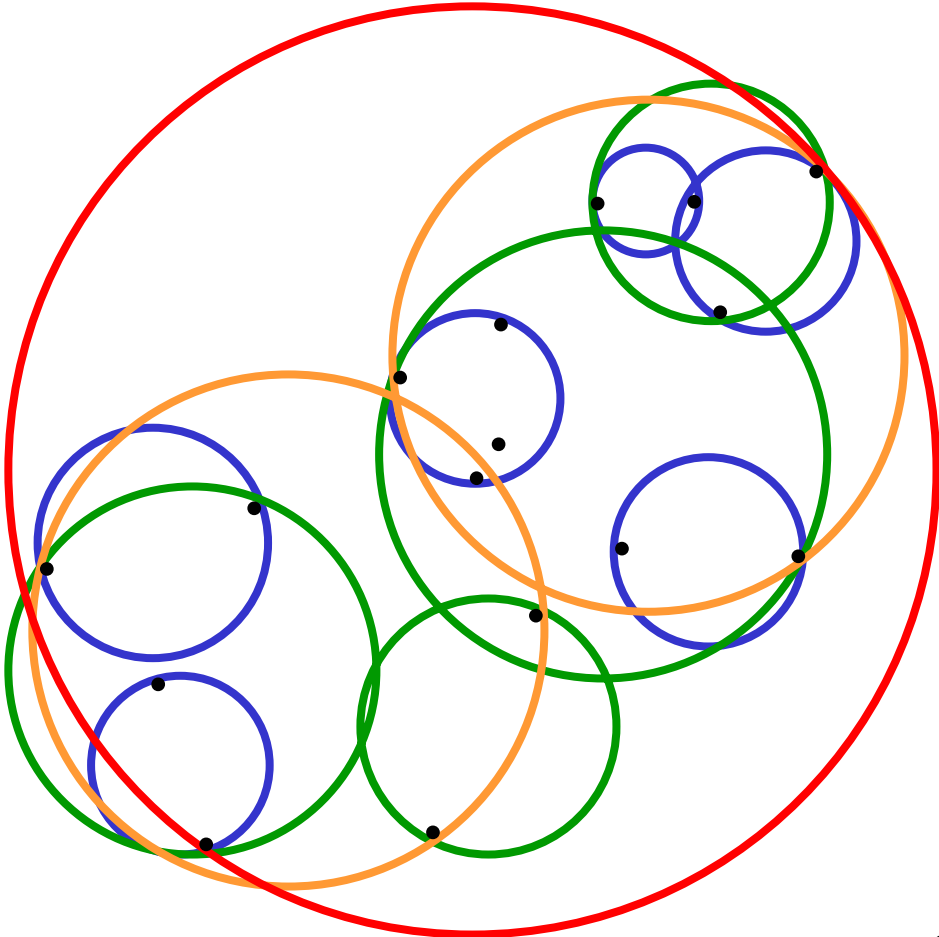
A Ball Tree



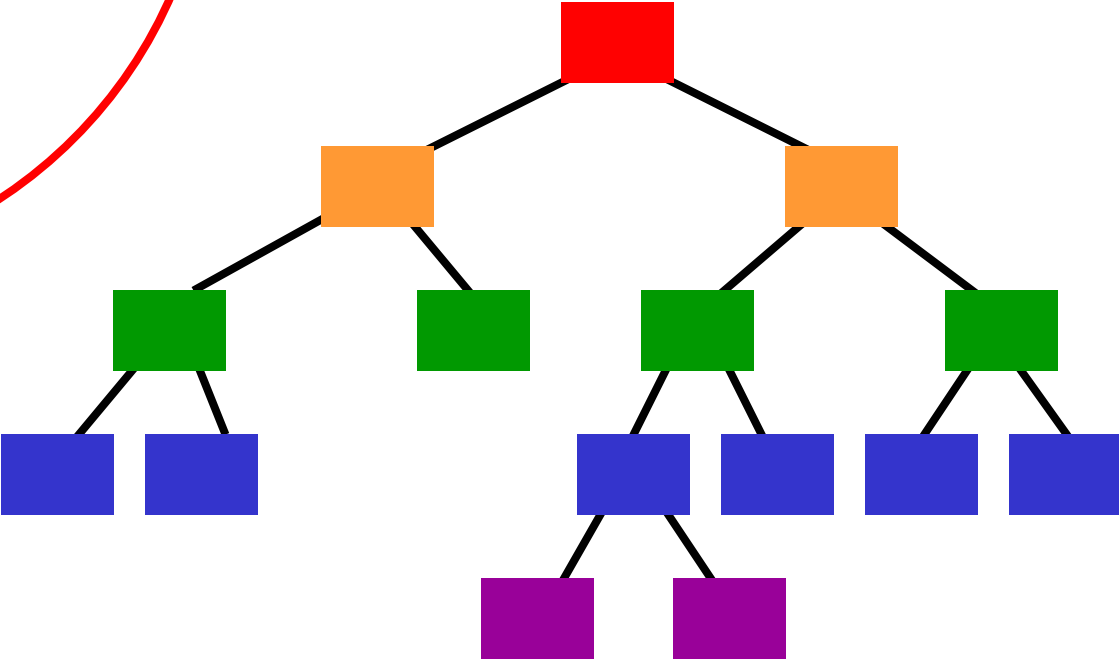
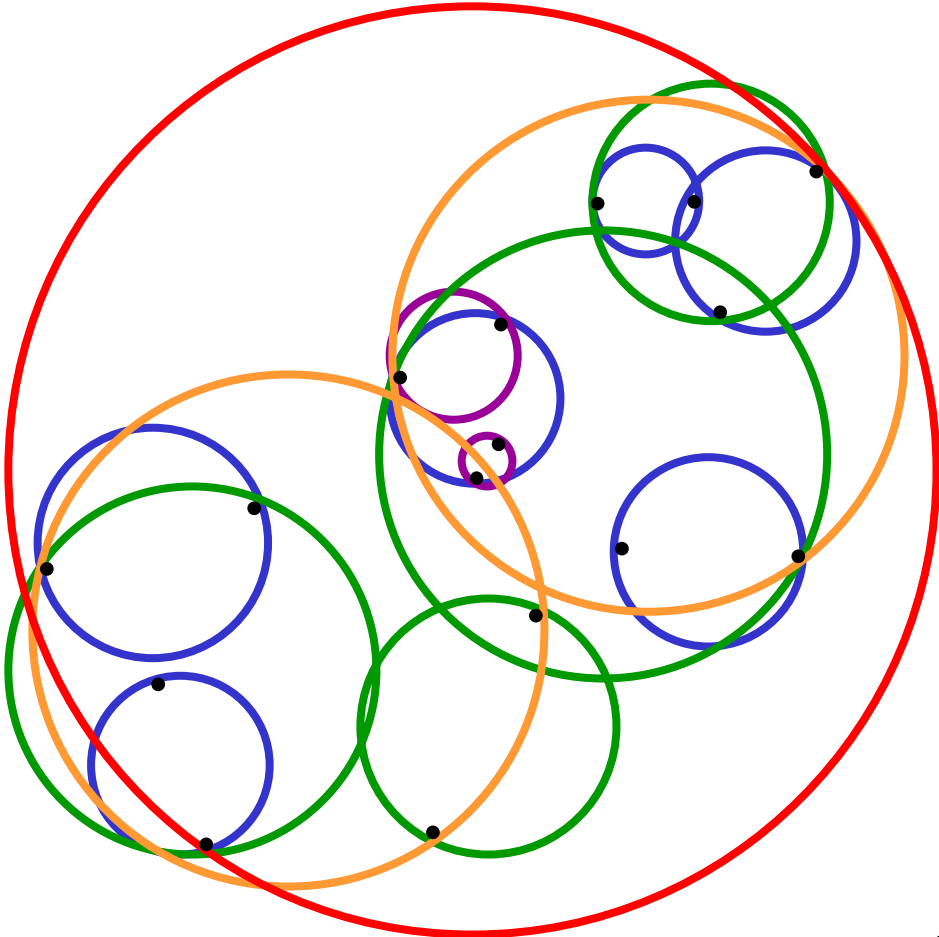
A Ball Tree



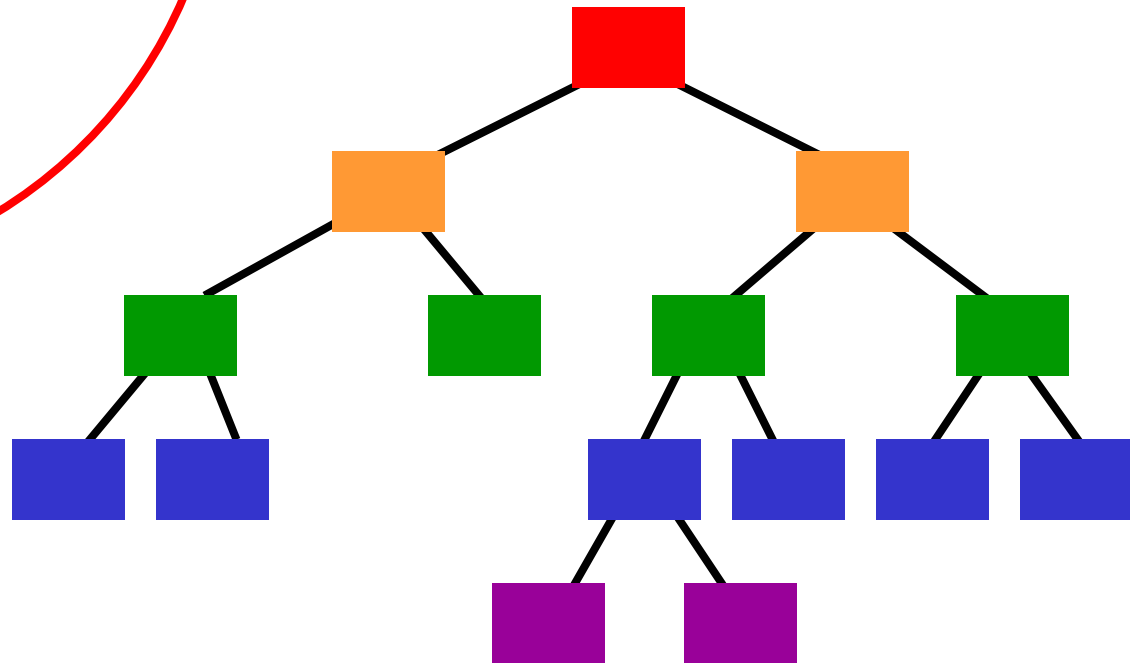
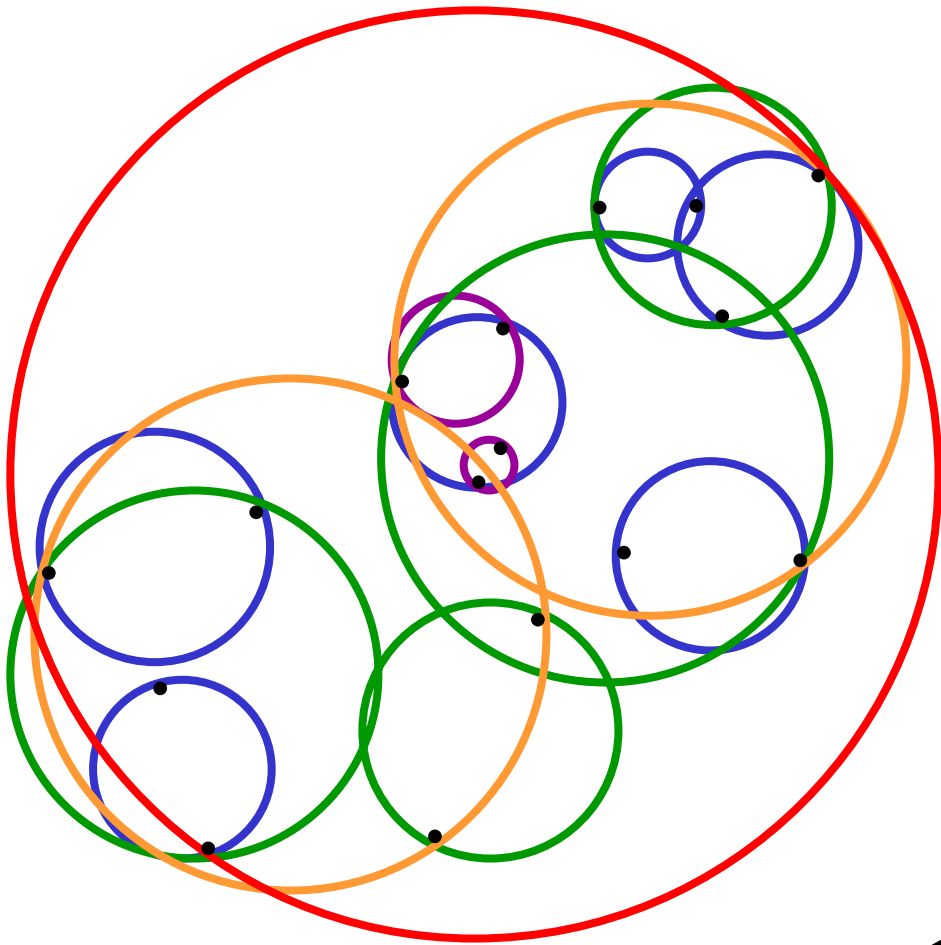
A Ball Tree



A Ball Tree



A Ball Tree



• J. Uhlmann, 1991

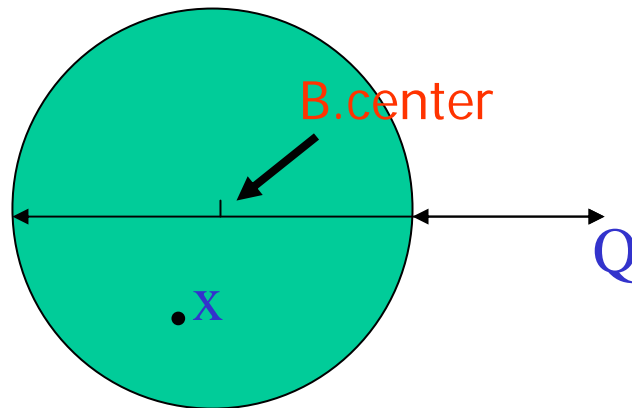
• S. Omohundro, NIPS 1991

Ball-trees: properties

Let Q be any query point and let x be a point inside ball B

$$|x - Q| \geq |Q - B.\text{center}| - B.\text{radius}$$

$$|x - Q| \leq |Q - B.\text{center}| + B.\text{radius}$$



Outline

Cached Sufficient Statistics

Ball Trees (= Metric Trees)

▶ K-nearest neighbor with ball trees

Very fast non-parametric classification

skewed binary outputs

General binary outputs

multi-classed outputs

Very fast kernel-based statistics

n-point computations

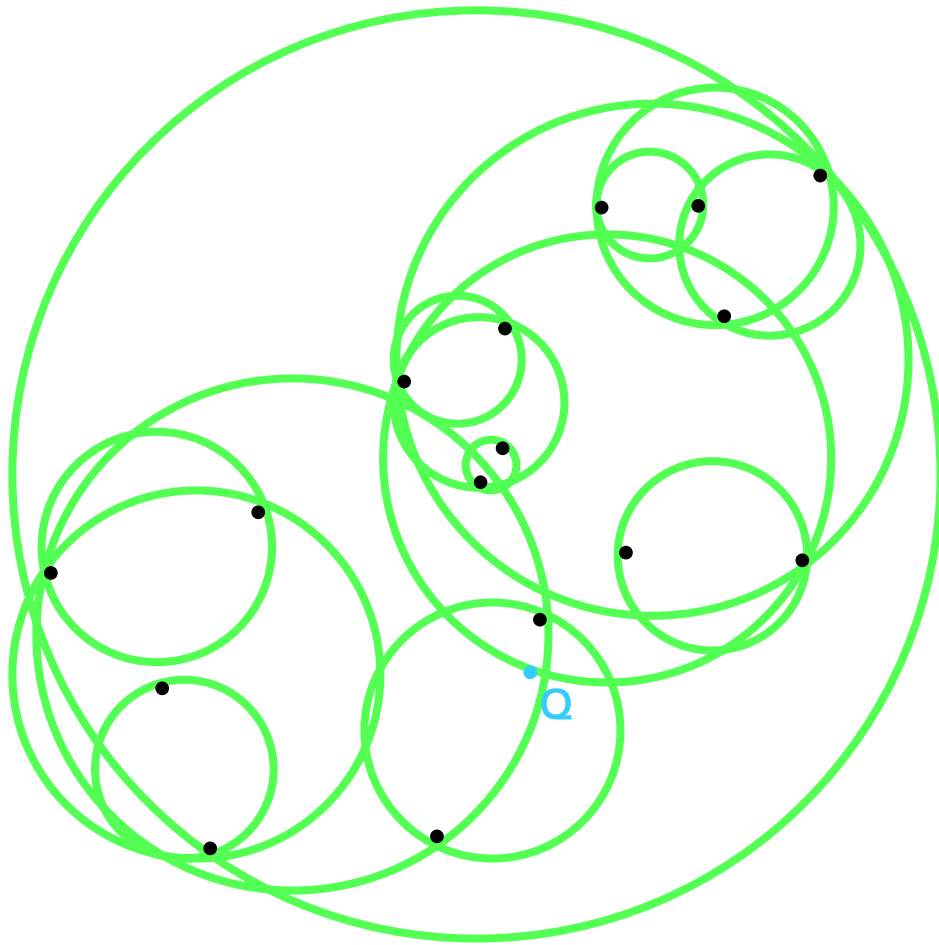
clustering

non-parametric clustering (overdensity hunting)

Active learning for anomaly hunting

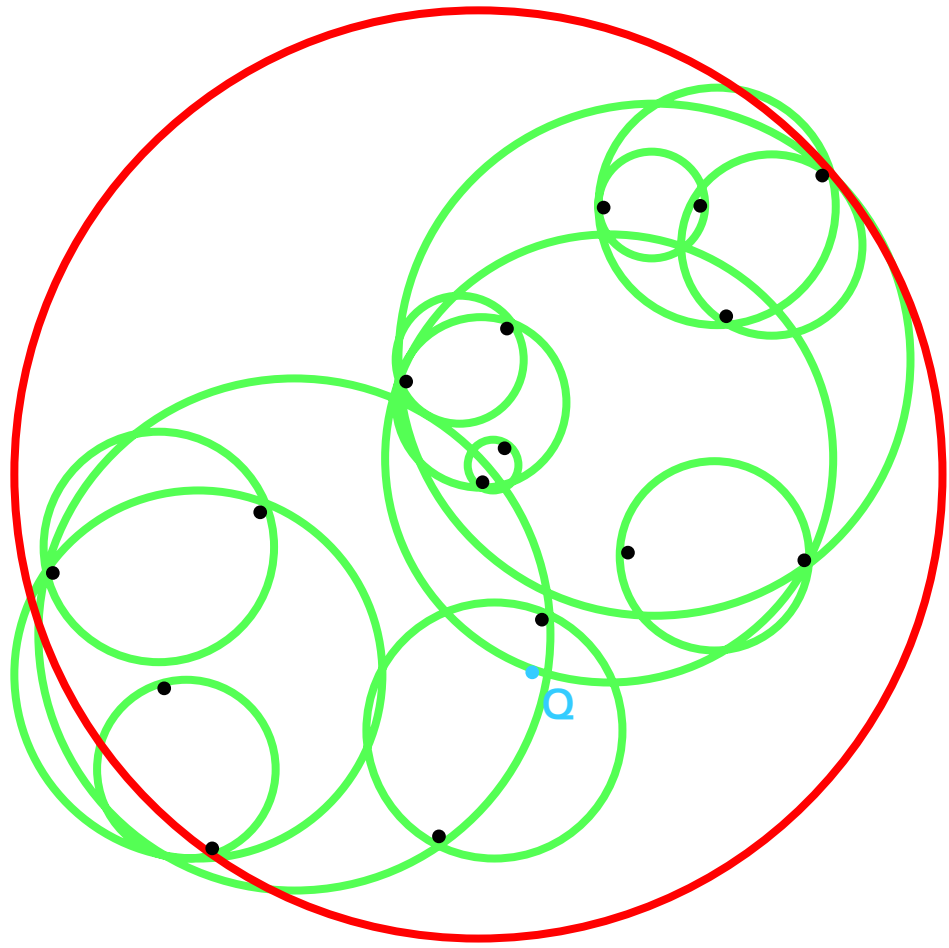
GMorph: Efficient Galaxy morphology fitting

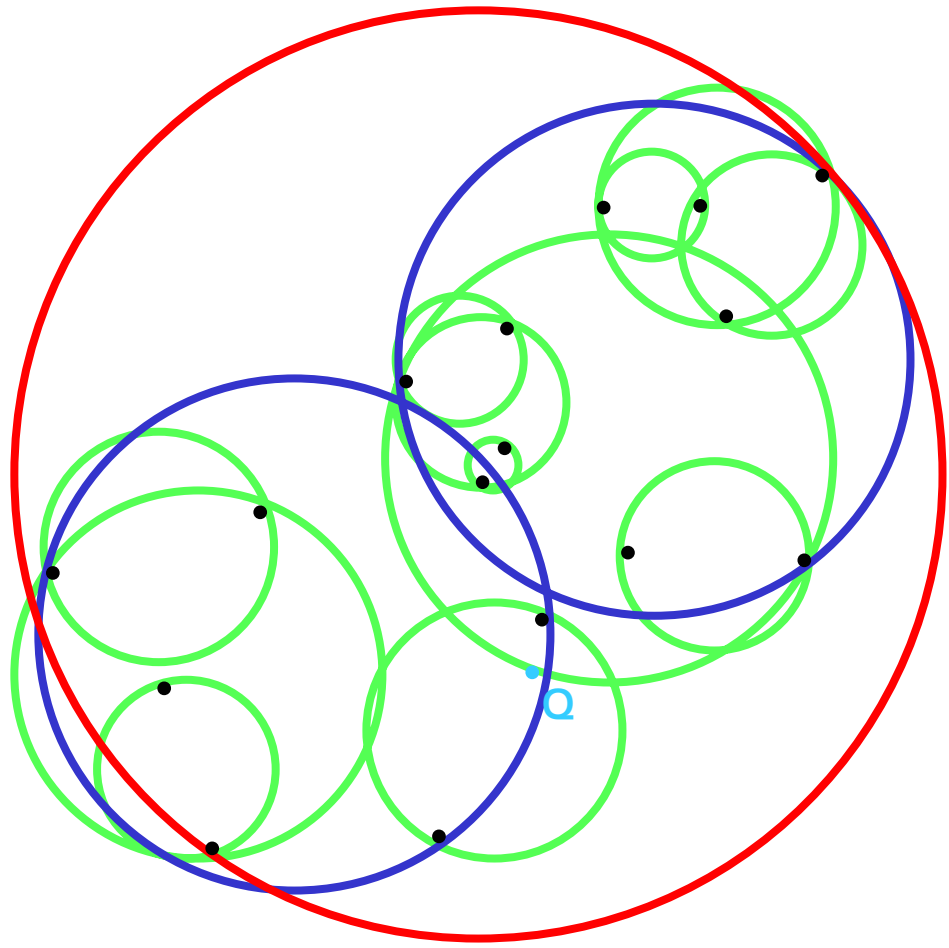
Other Auton topics

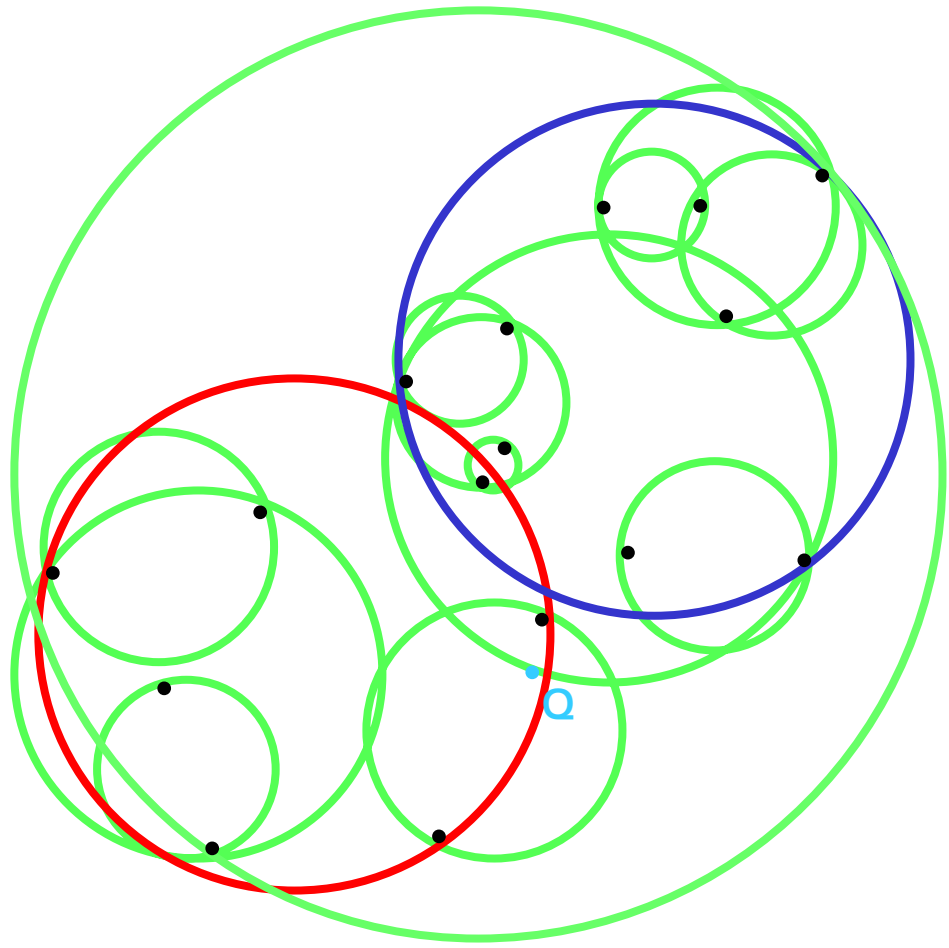


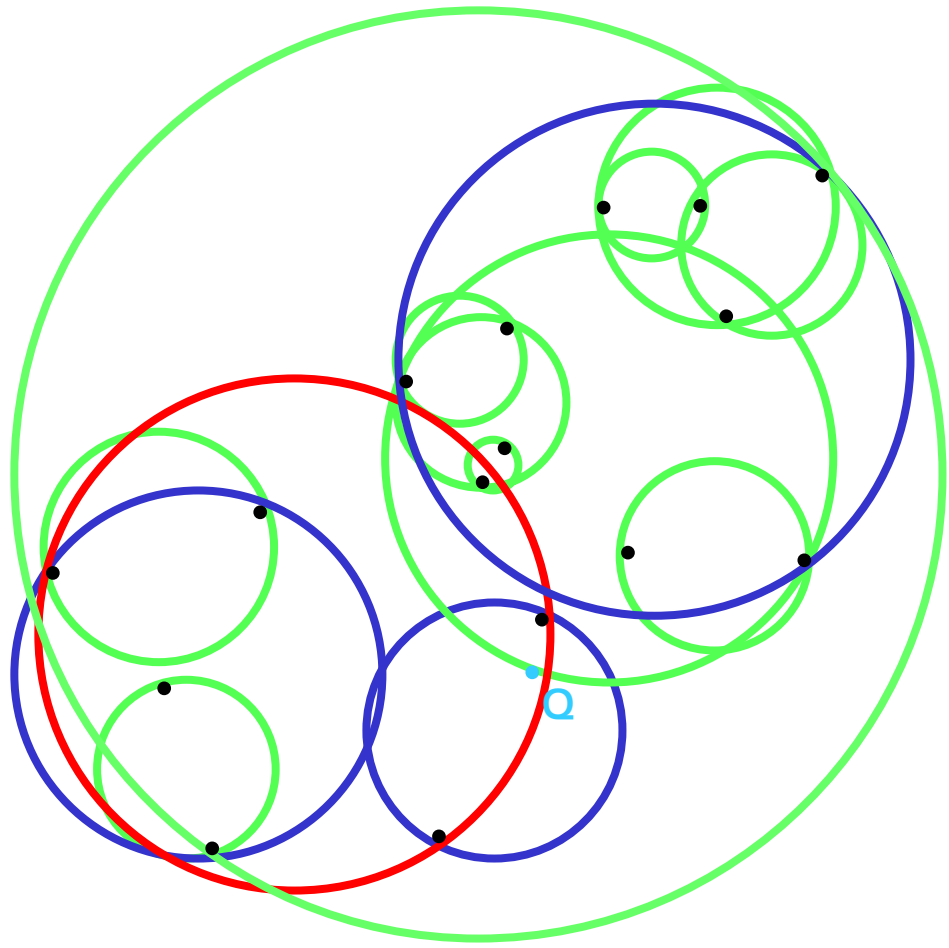
Goal: Find out
the 2-nearest
neighbors of Q .

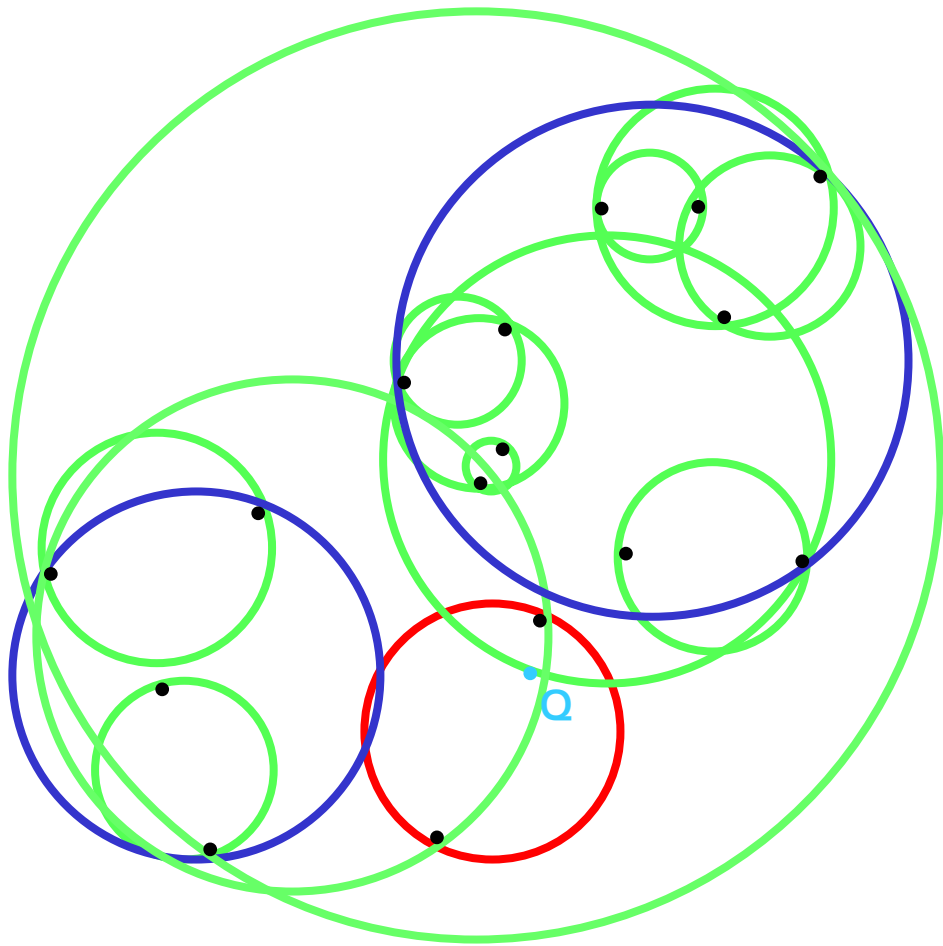
- J. Uhlmann, 1991
- S. Omohundro, NIPS 1991

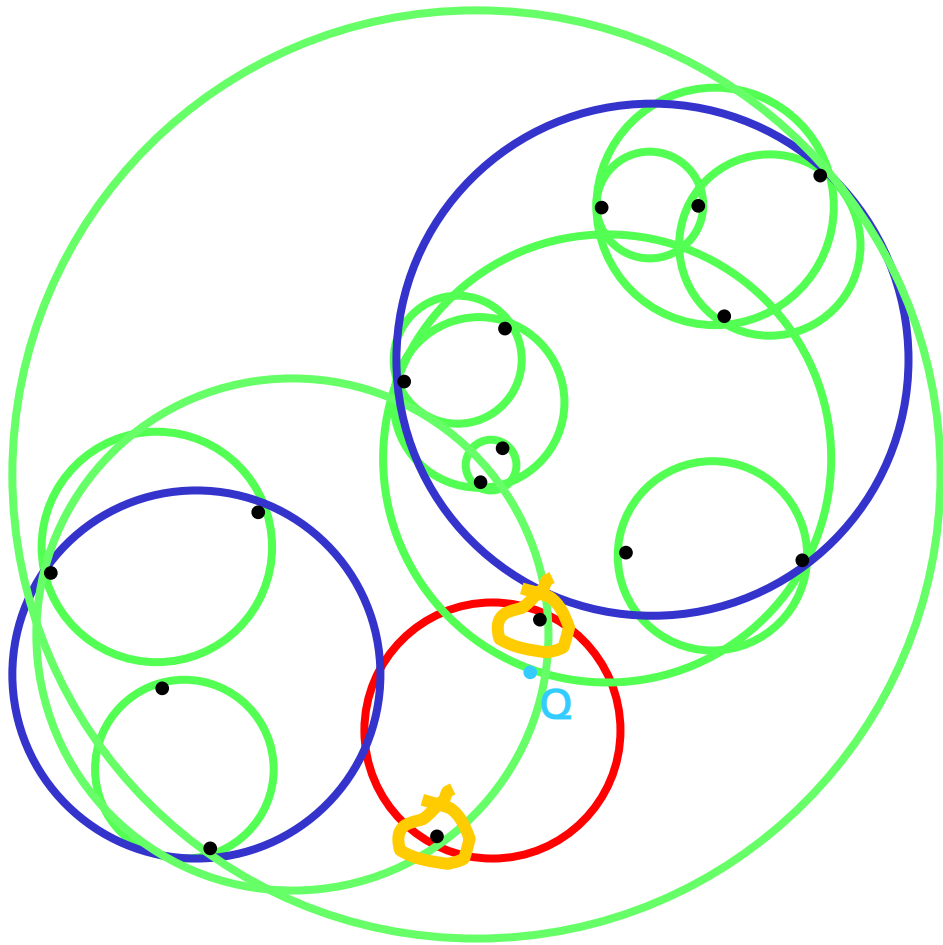




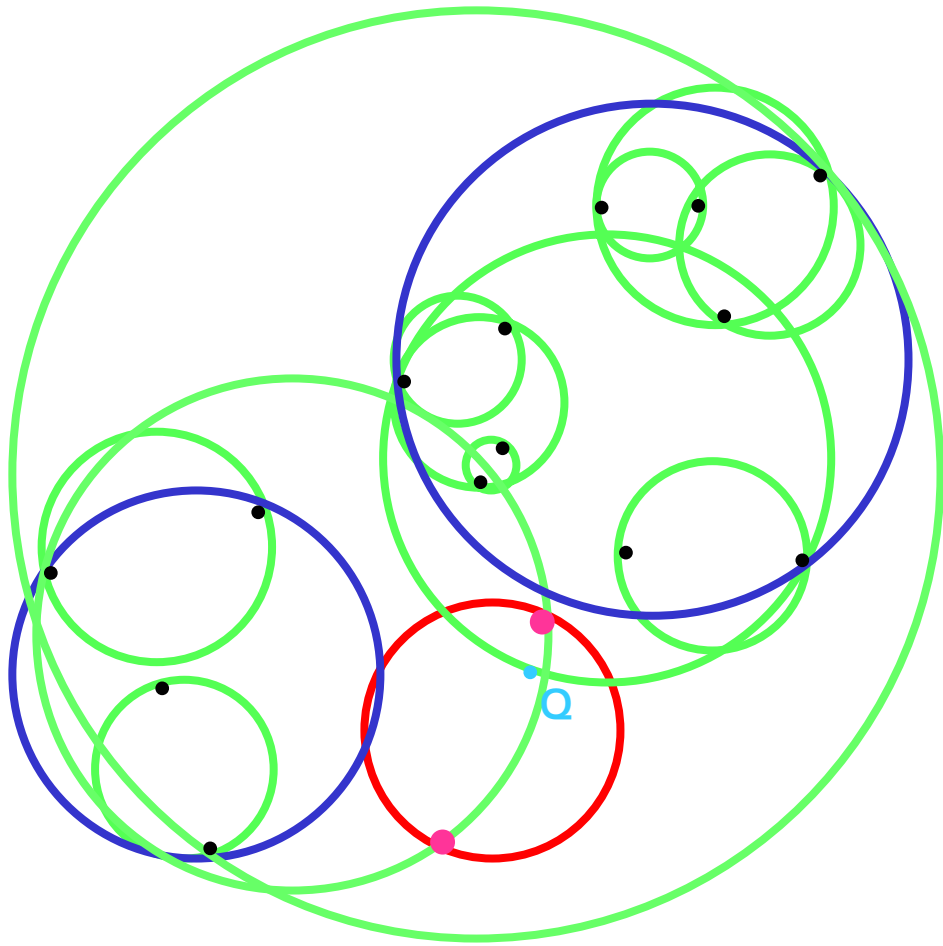






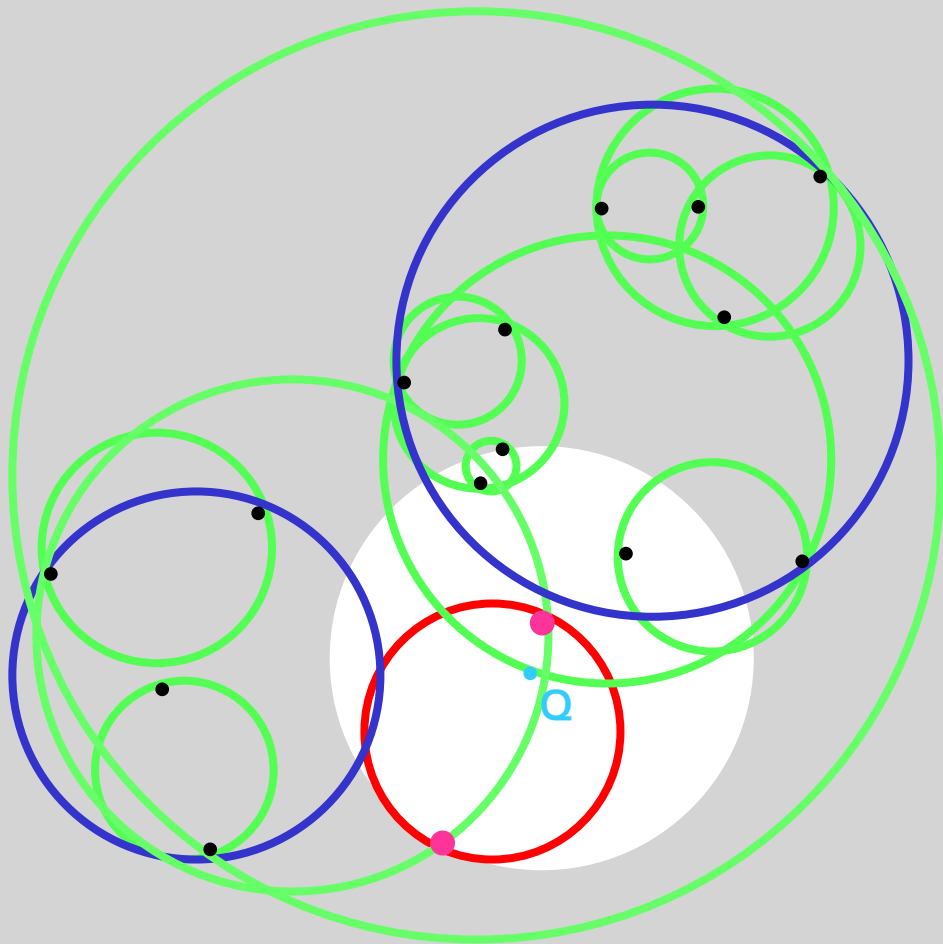


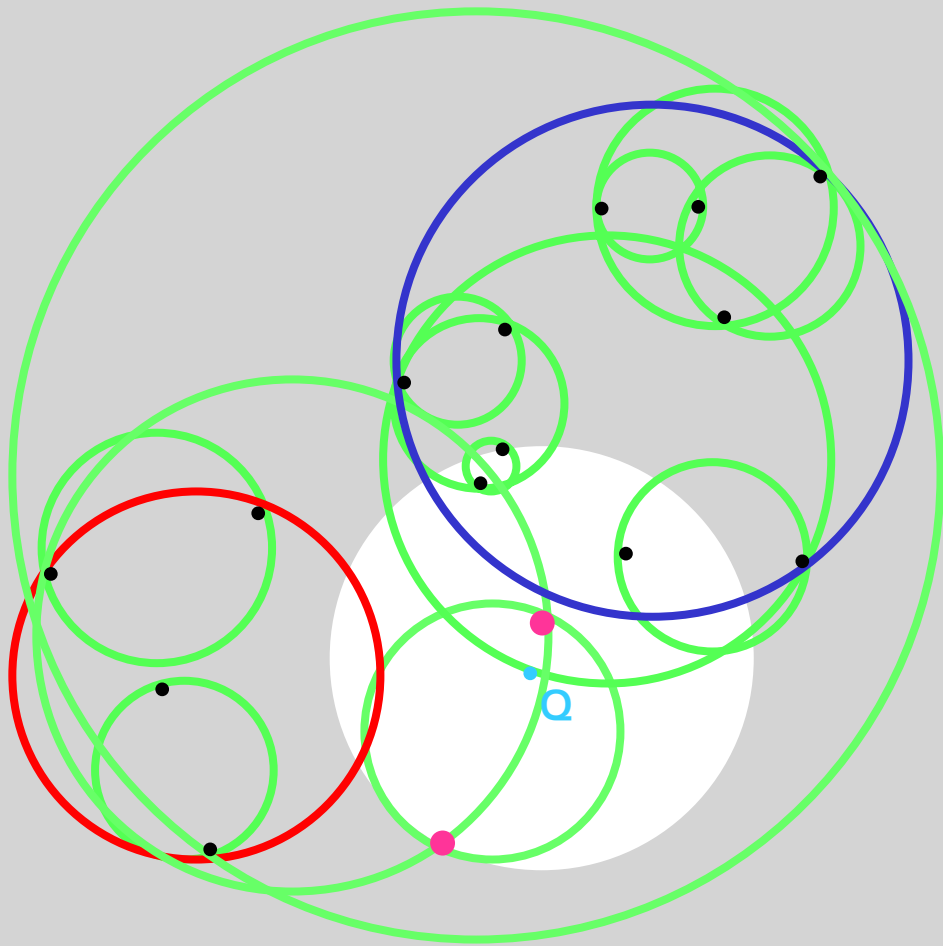
We've hit a leaf node, so we explicitly look at the points in the node

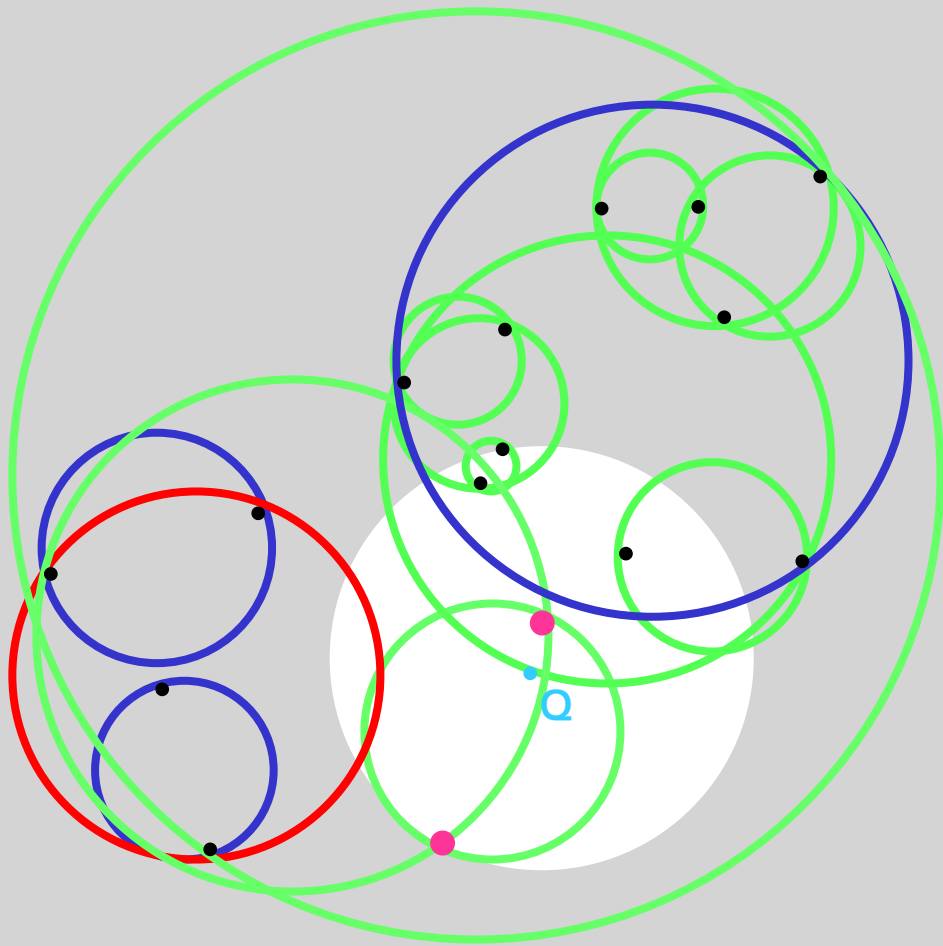


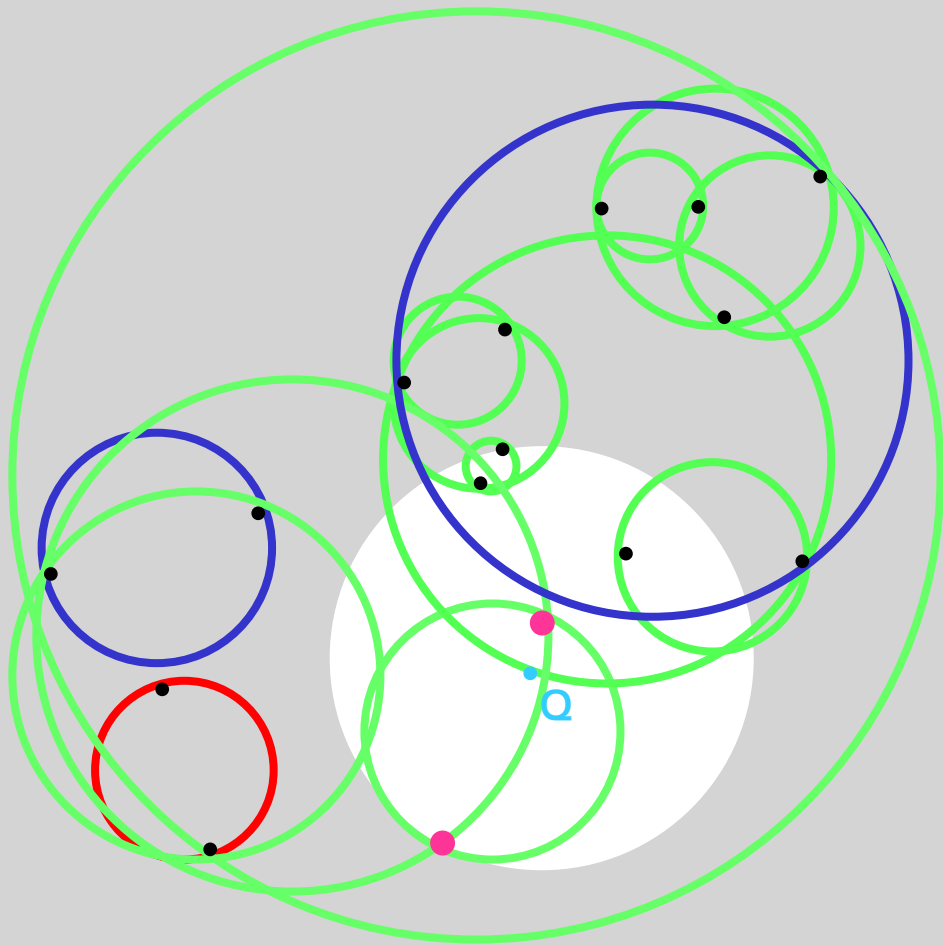
Two nearest
neighbors found so
far are in pink

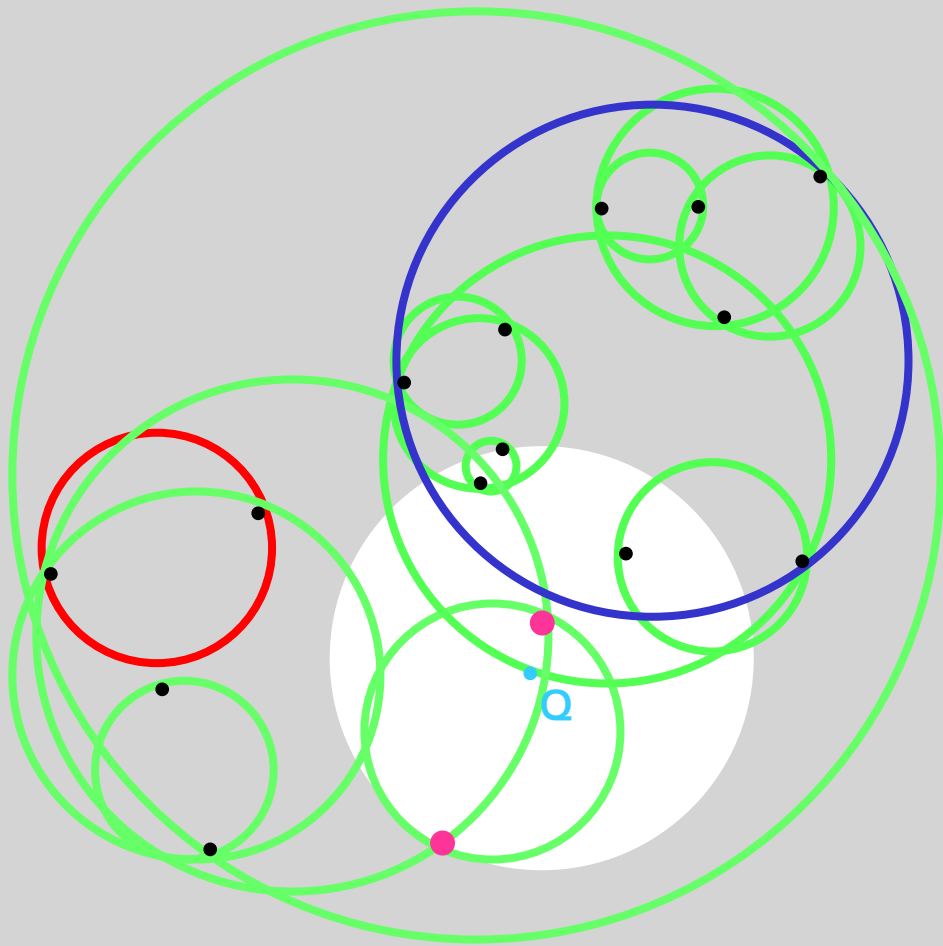
(remember we have
yet to search the
blue balls)

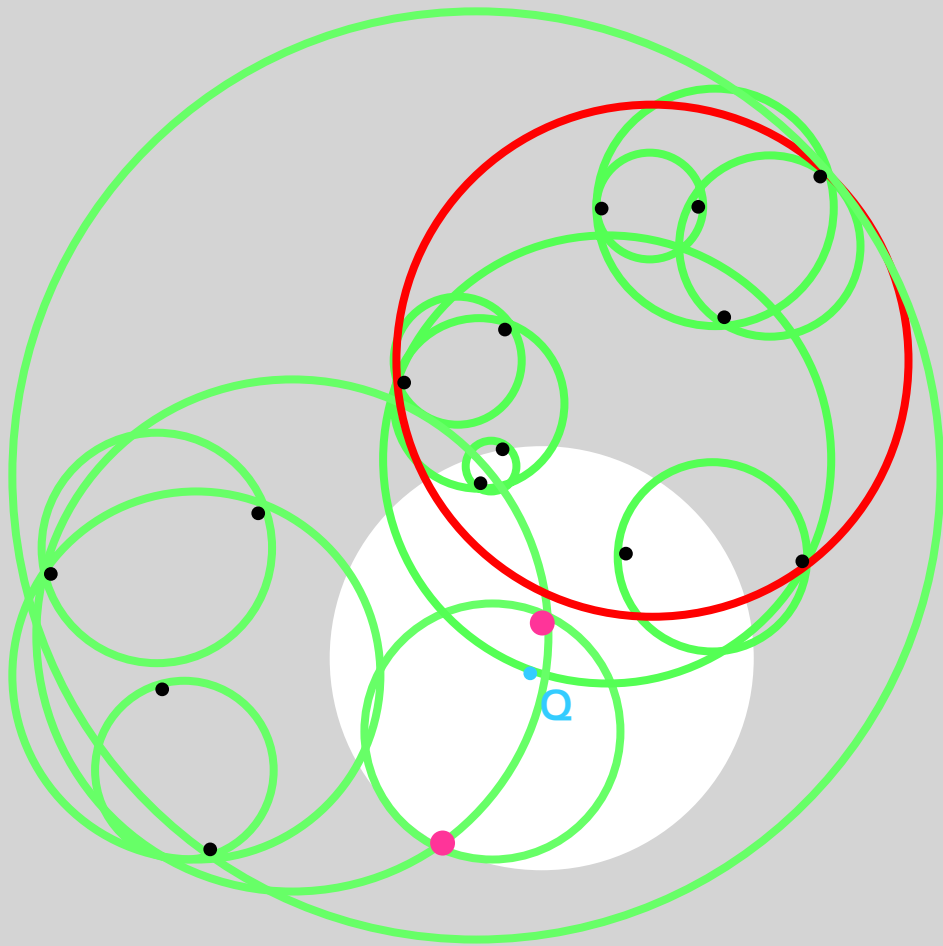


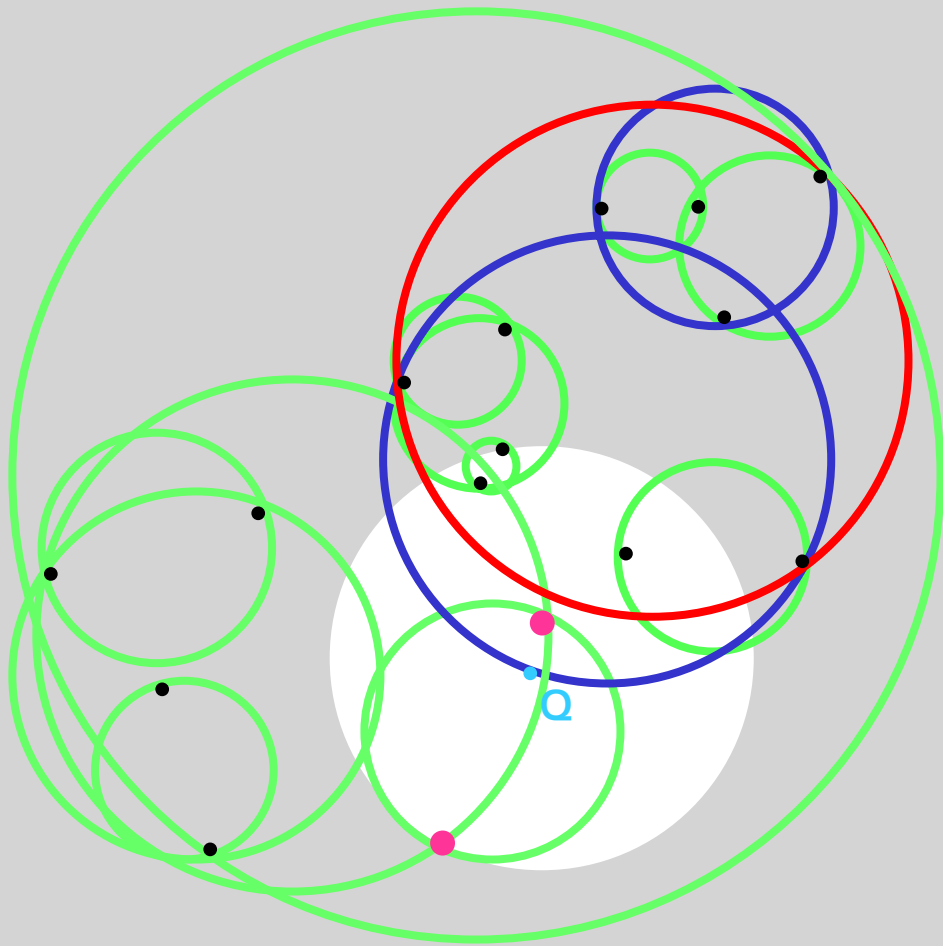


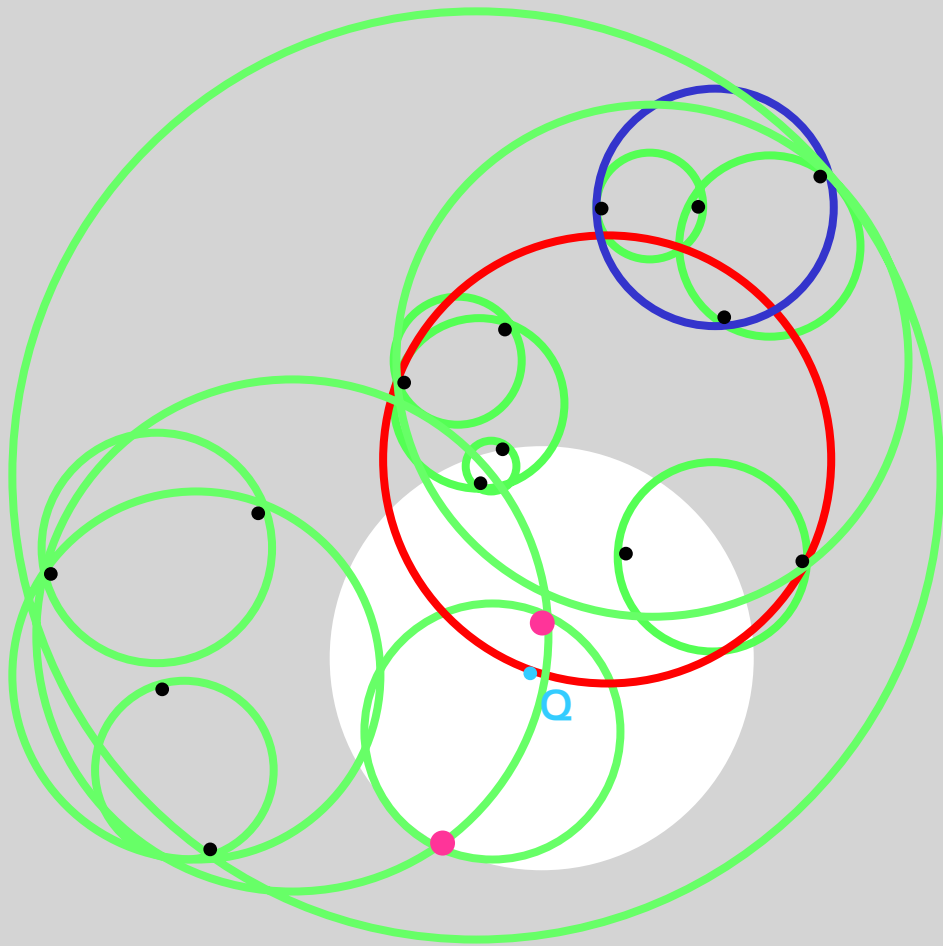


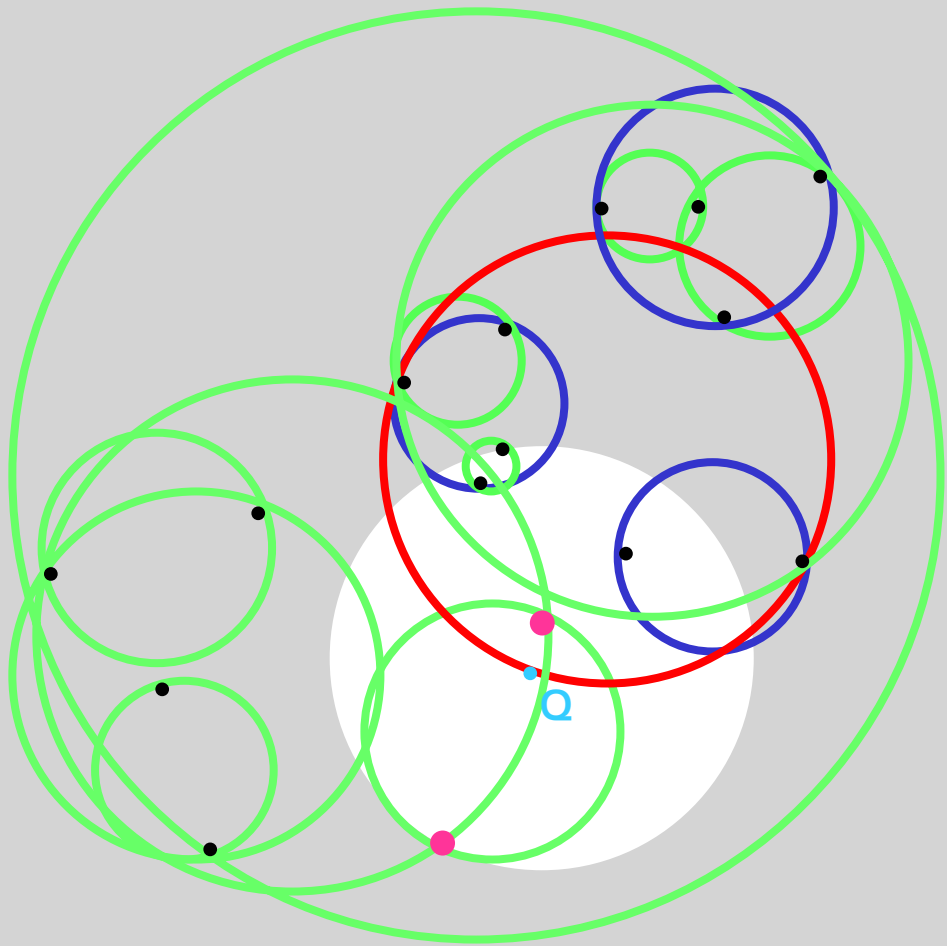


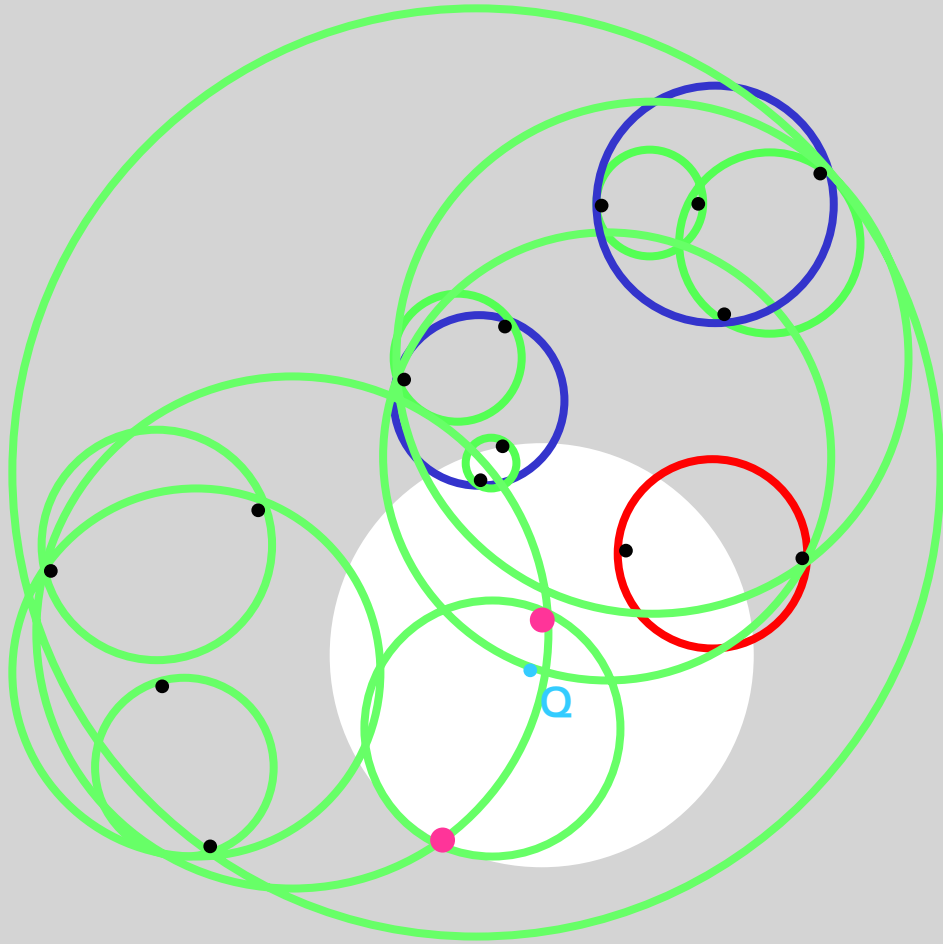


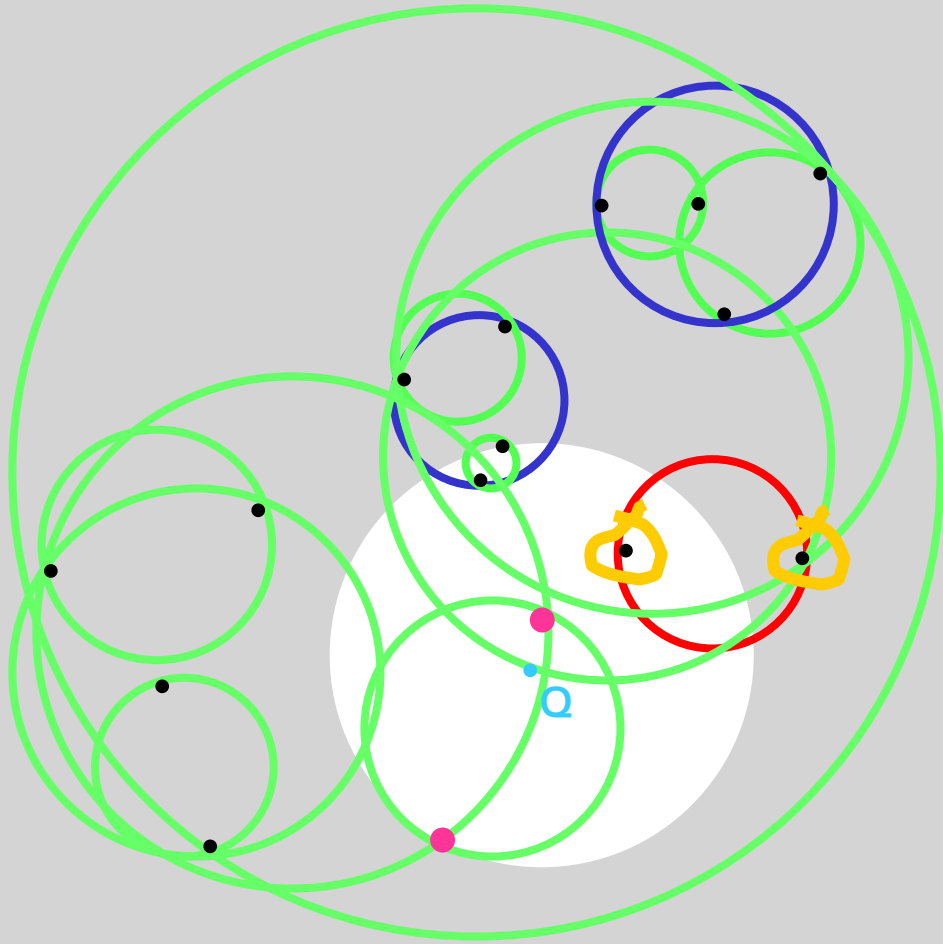


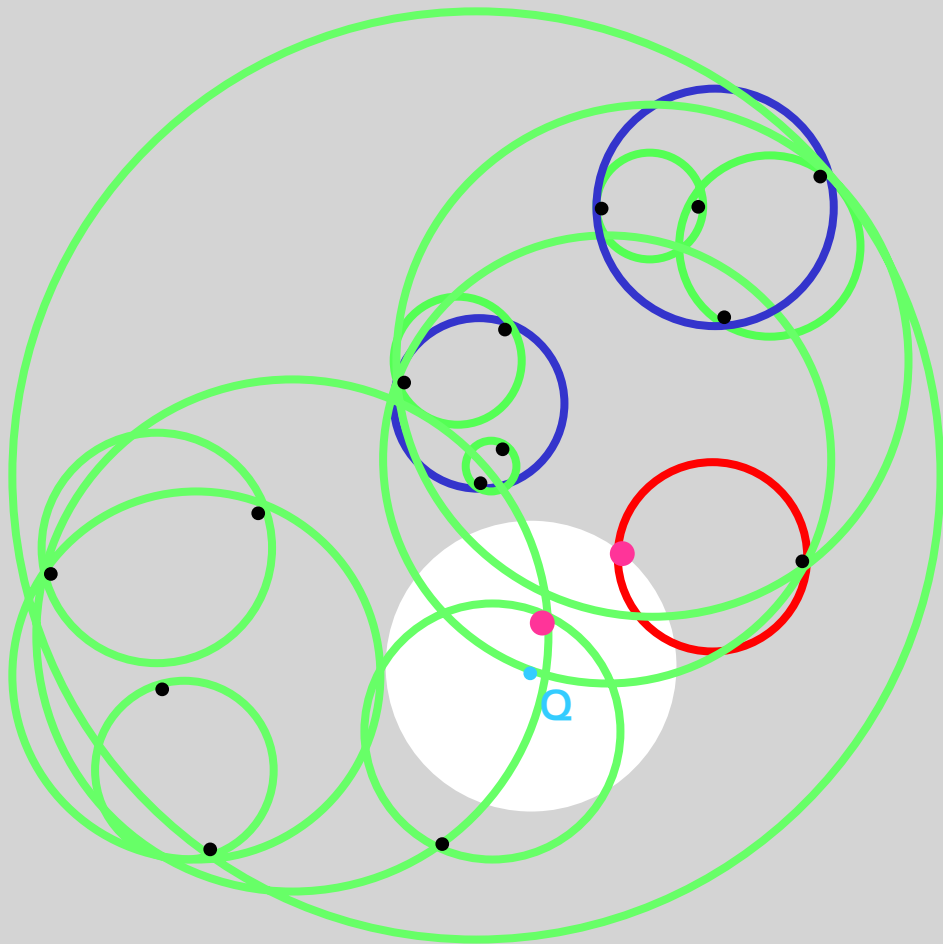


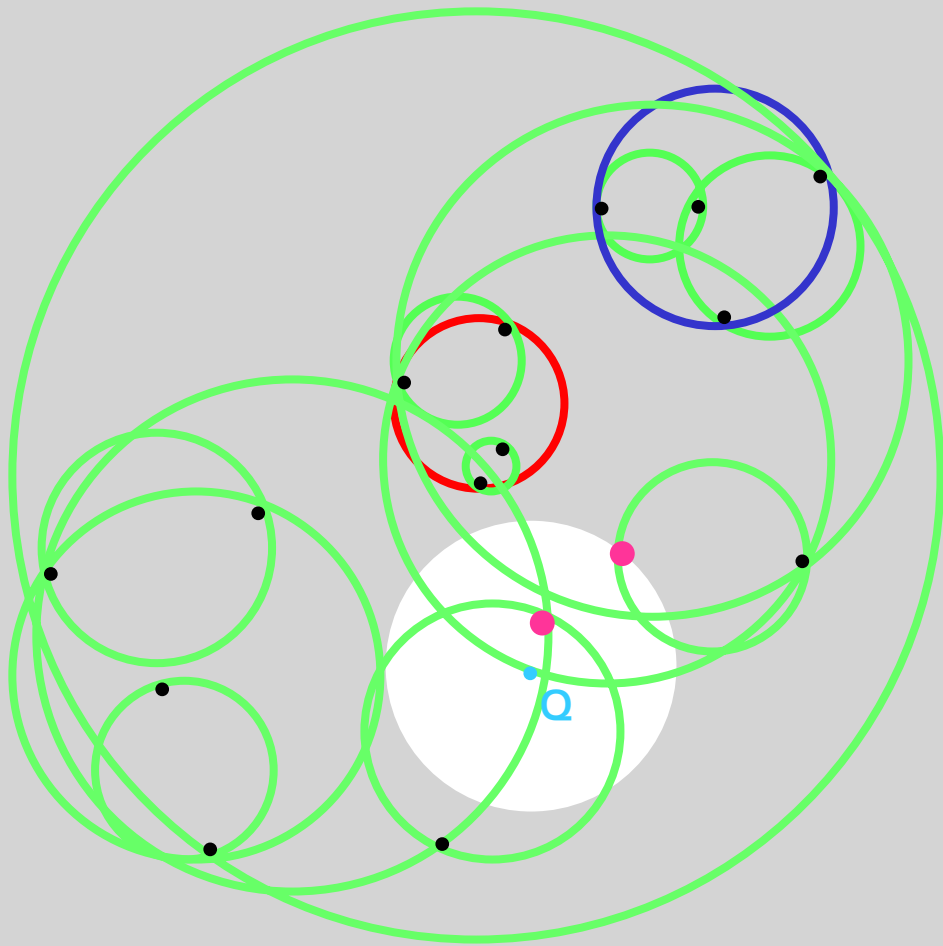


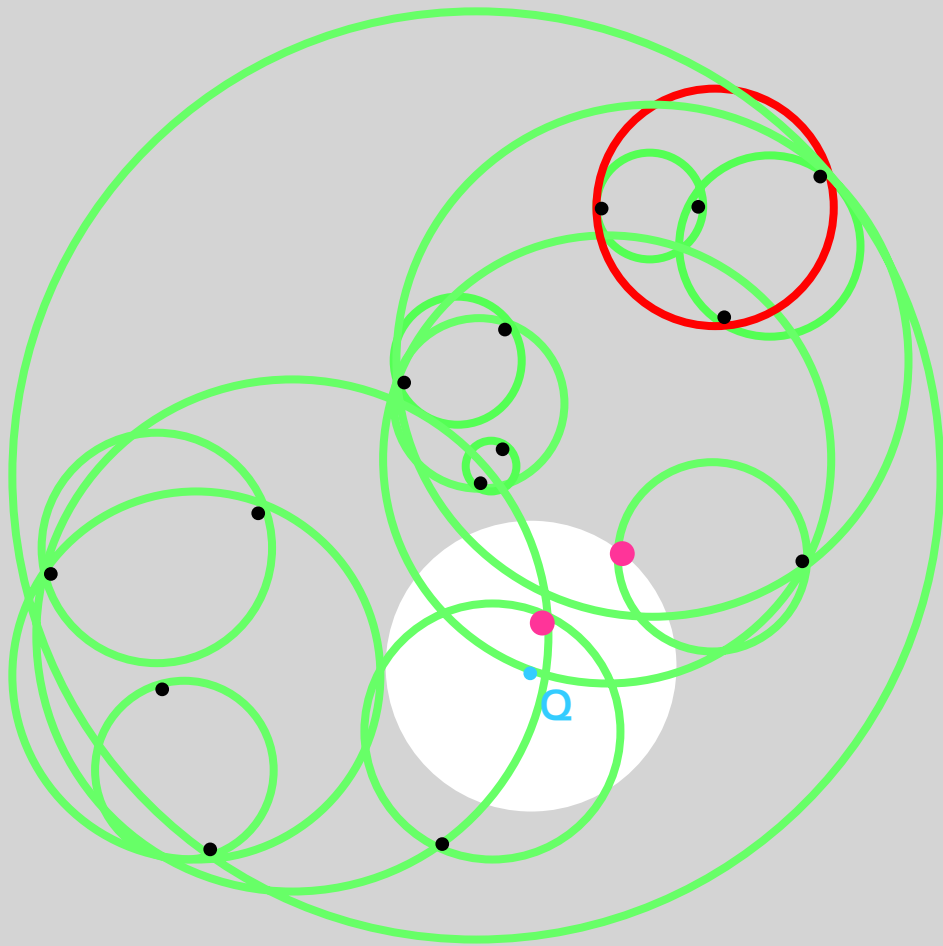


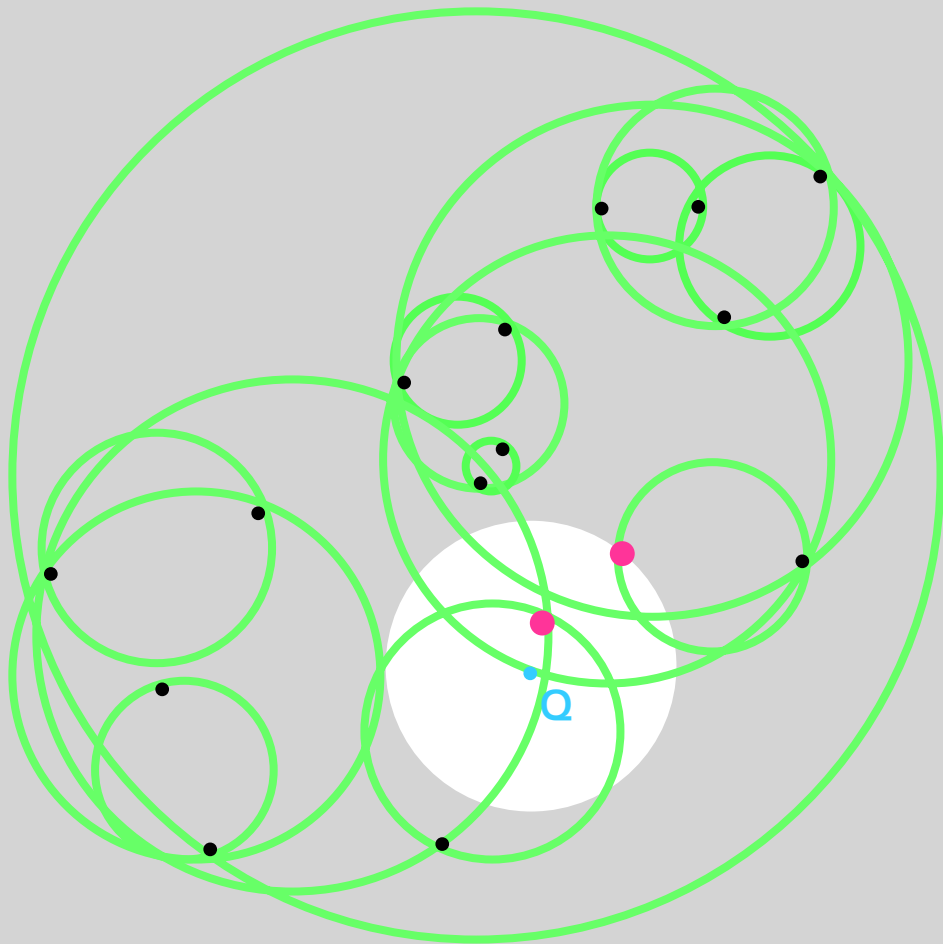












Outline

Cached Sufficient Statistics

Ball Trees (= Metric Trees)

K-nearest neighbor with ball trees



Very fast non-parametric classification

skewed binary outputs

General binary outputs

multi-classed outputs

Very fast kernel-based statistics

n-point computations

clustering

non-parametric clustering (overdensity hunting)

Active learning for anomaly hunting

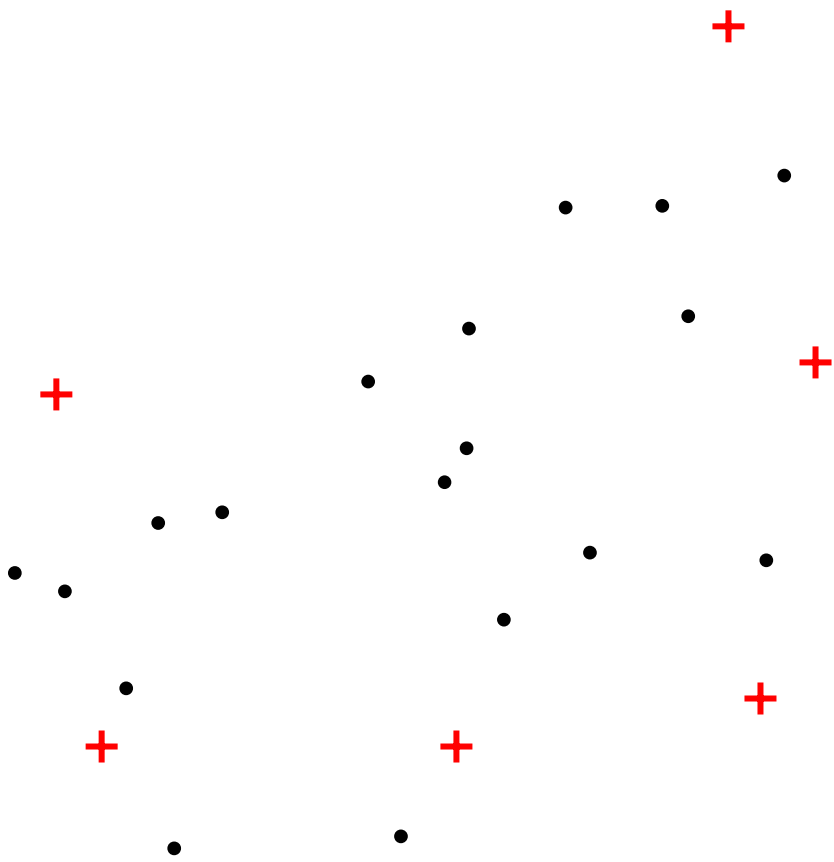
GMorph: Efficient Galaxy morphology fitting

Other Auton topics

KNS2

- Assume binary output
- Assume positive class is much less frequent than negative class
- Assume we want more than a “positive/negative” prediction: we want to know exactly how many of the K-NN are from the +ve class

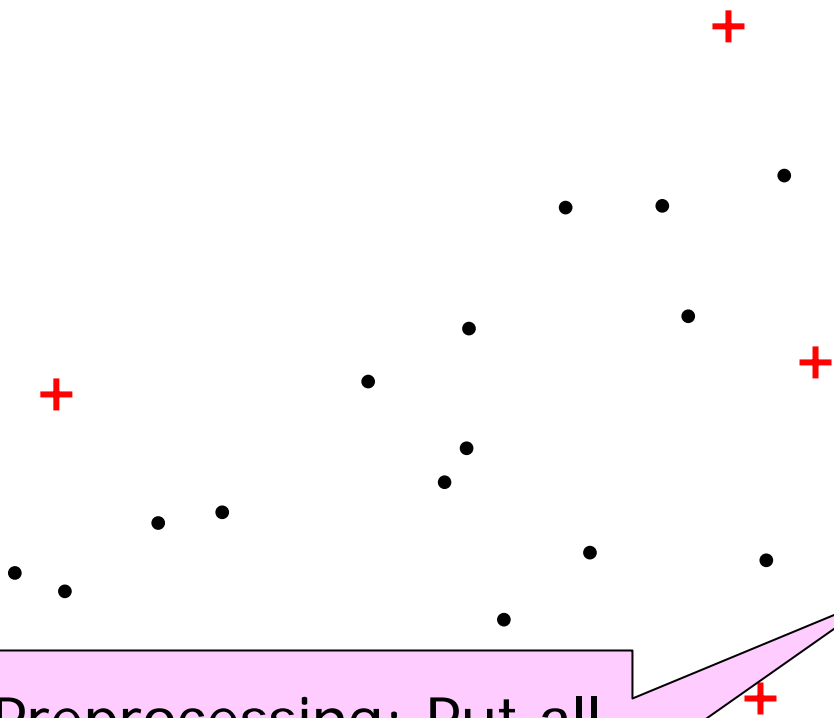
KNS2 does this without finding the K-NN



Assume we have a set of data points.

Some are +ve points (denoted **+**)

The large majority are -ve points (denoted **●**)



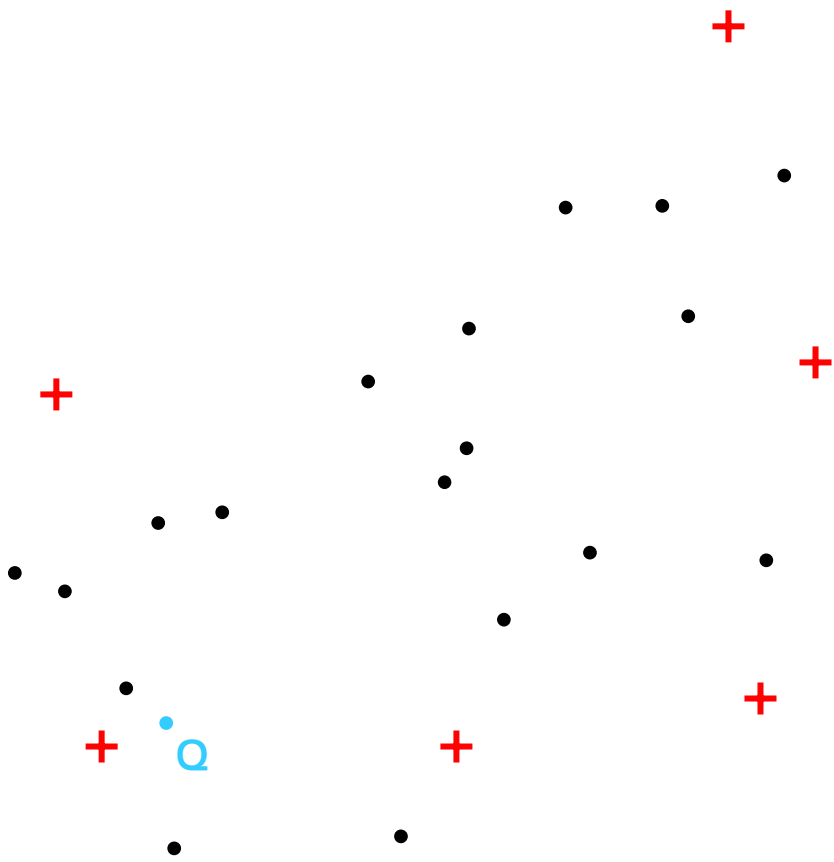
Assume we have a set of data points.

Some are +ve points (denoted **+**)

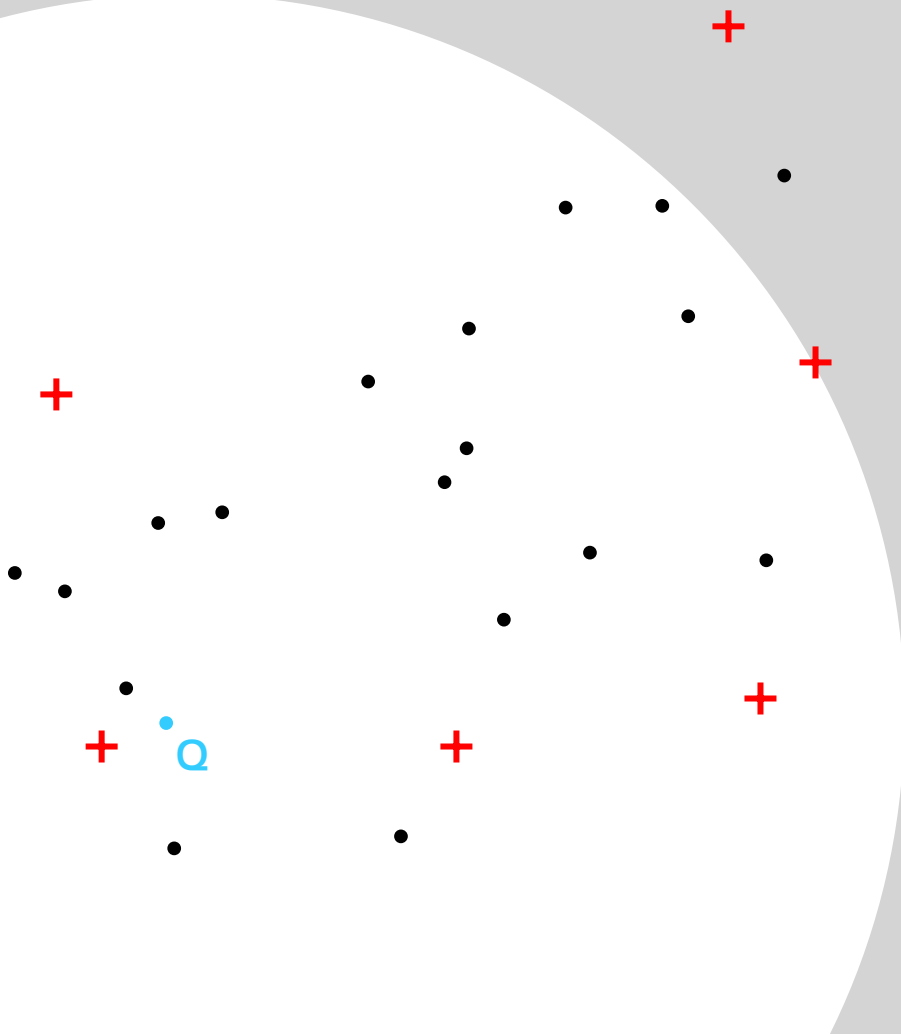
The vast majority are -ve points (denoted **•**)

Preprocessing: Put all your +ve points in a small ball tree

Preprocessing: Put all your -ve points in a separate large ball tree

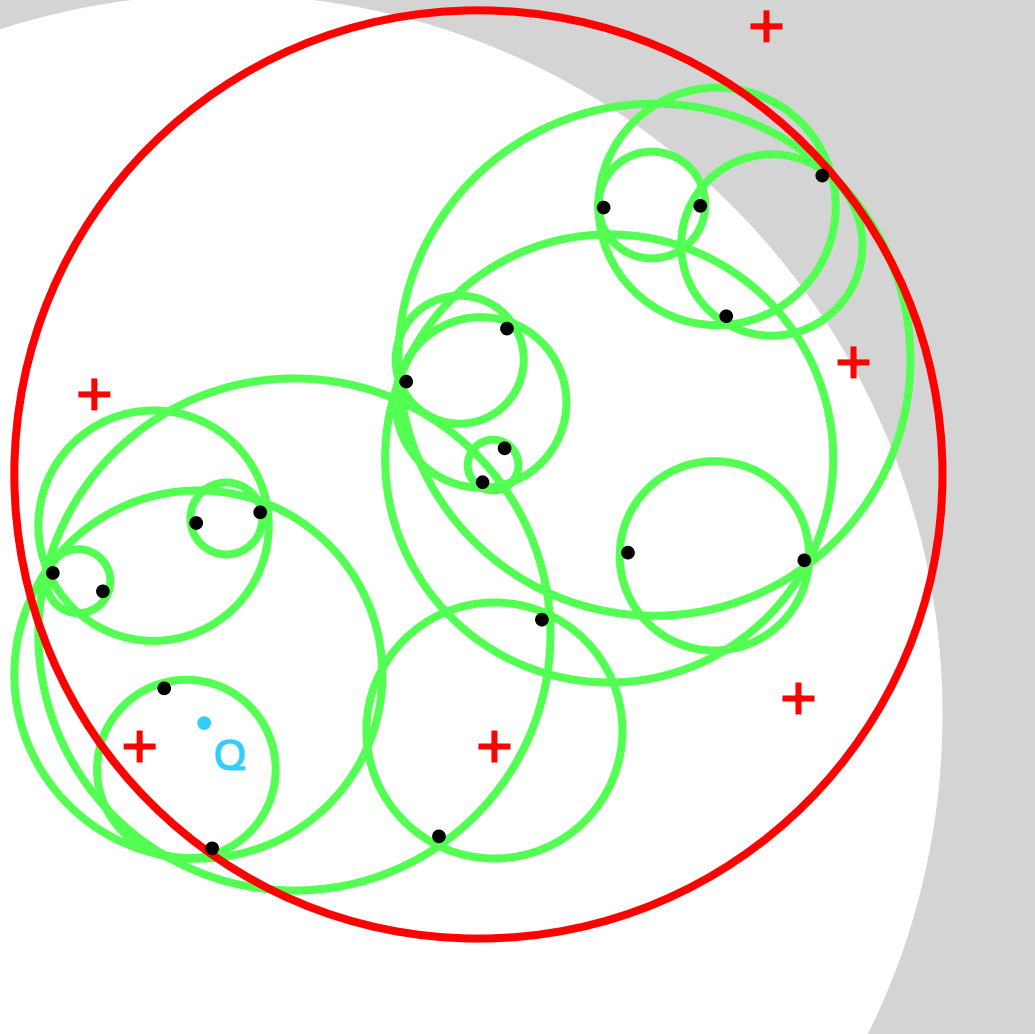


Goal: Find out how many of the 5-nearest neighbors of Q are positive.

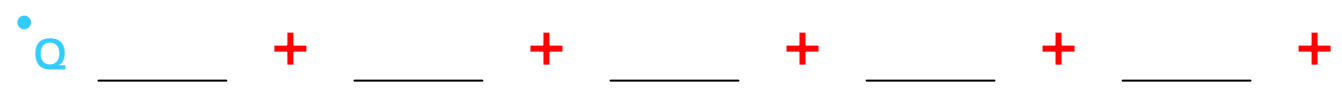


Step One: Find the five nearest +ve points using KNS1.

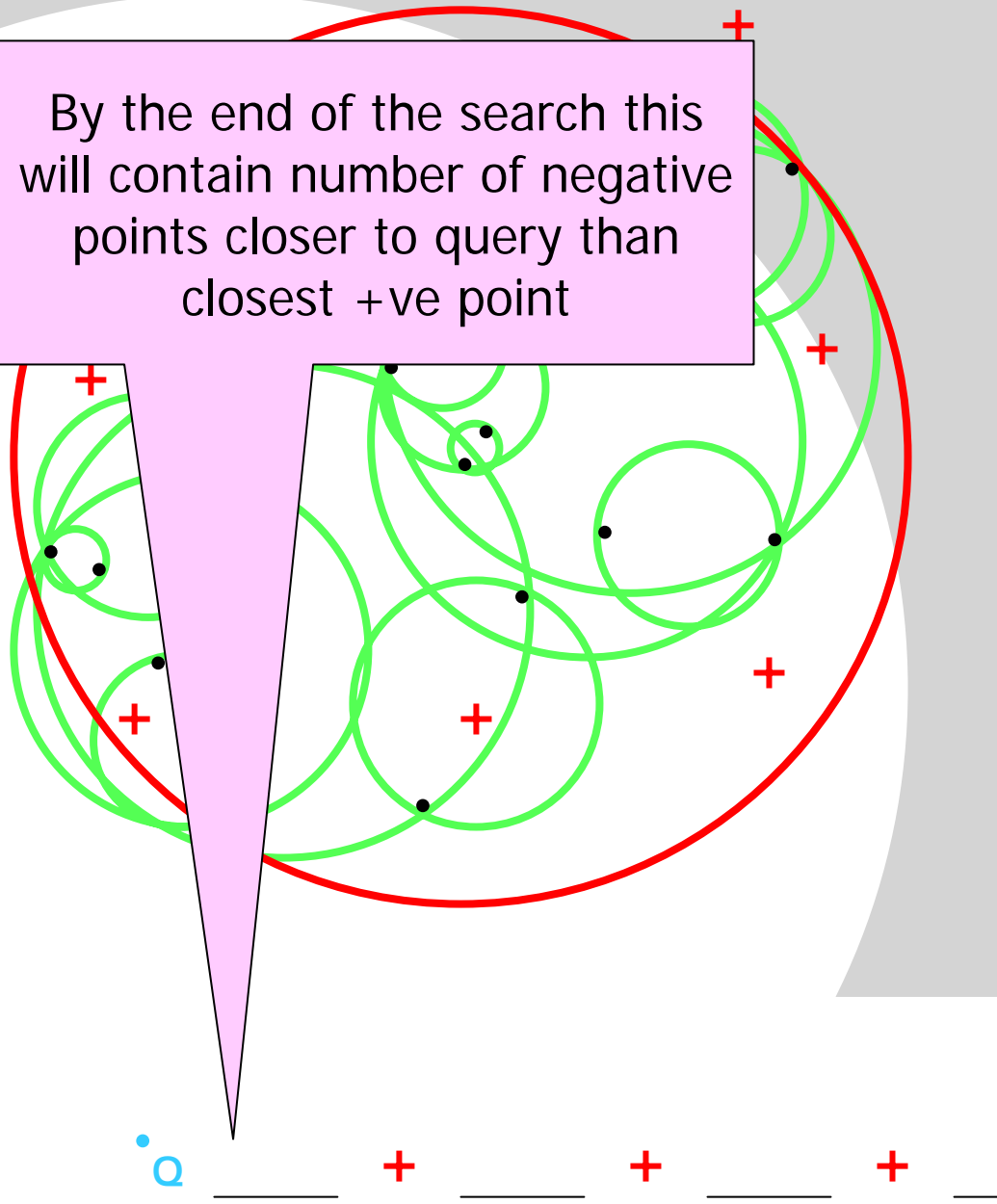
We're assuming there are far fewer +ves than -ves so this is not the dominant cost.



Step 2: Search the ball-tree of -ve points starting at the root.



By the end of the search this will contain number of negative points closer to query than closest +ve point



Search the ball-tree of -ve points starting at the root.

By the end of the search this will contain number of negative points closer to query than closest +ve point

By the end of the search this will contain number of negative points whose distance to query is between distances of closest +ve point and 2nd closest +ve point

Search the ball-tree of -ve points starting the root.

q

+

+

+

+

+

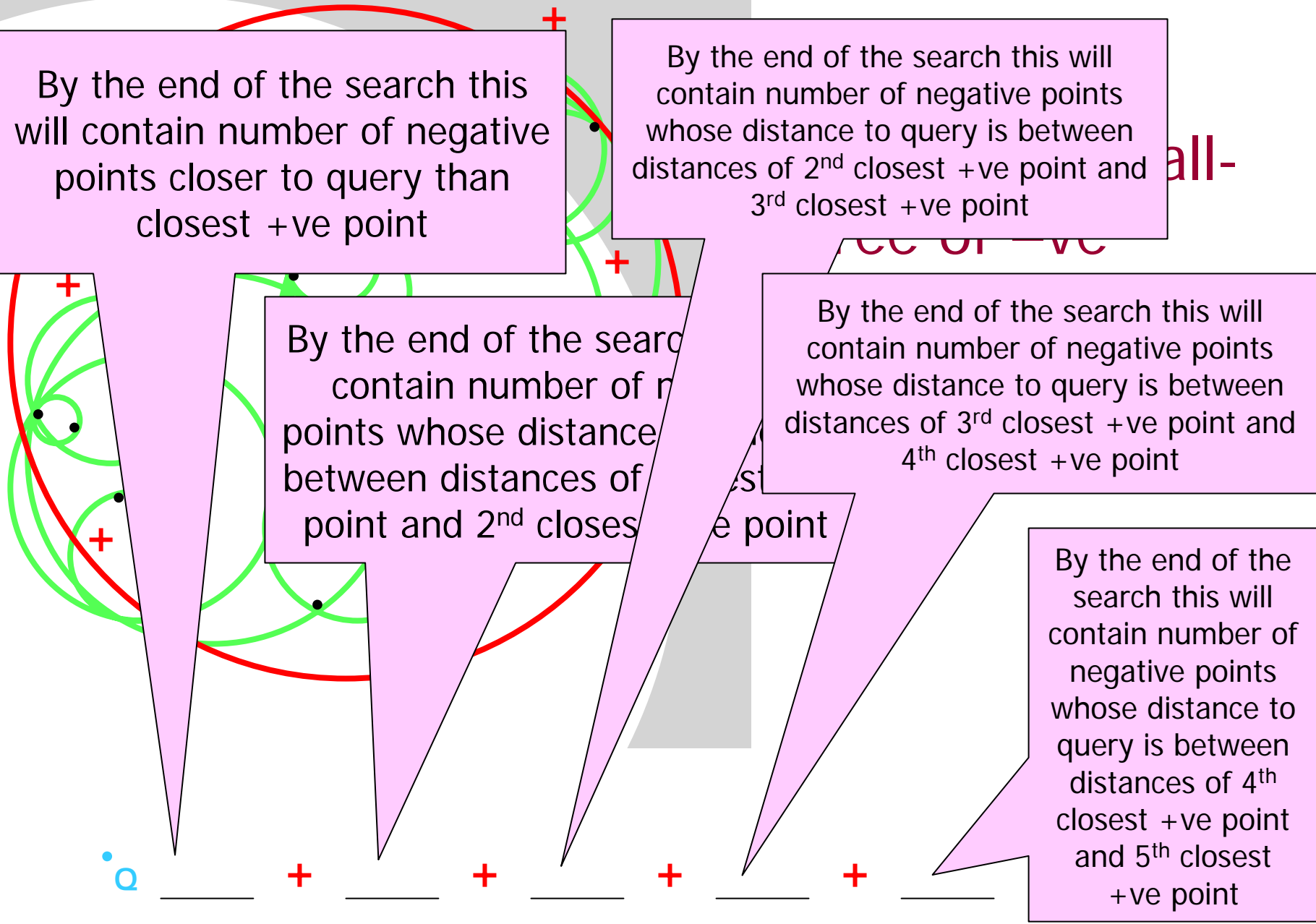
By the end of the search this will contain number of negative points closer to query than closest +ve point

By the end of the search this will contain number of negative points whose distance to query is between distances of 2nd closest +ve point and 3rd closest +ve point

By the end of the search contain number of negative points whose distance between distances of point and 2nd closest +ve point

By the end of the search this will contain number of negative points whose distance to query is between distances of 3rd closest +ve point and 4th closest +ve point

By the end of the search this will contain number of negative points whose distance to query is between distances of 4th closest +ve point and 5th closest +ve point



By the end of the search this will contain number of negative points closer to query than the closest

But only if relevant to 5-NN query!

By the end of the search this will contain number of negative points whose distance to query is between distances of 2nd closest +ve point and 3rd closest +ve point

But only if relevant to 5-NN query!

By the end of the search this will contain number of negative points whose distance to query is between distances of 3rd and 4th closest +ve point

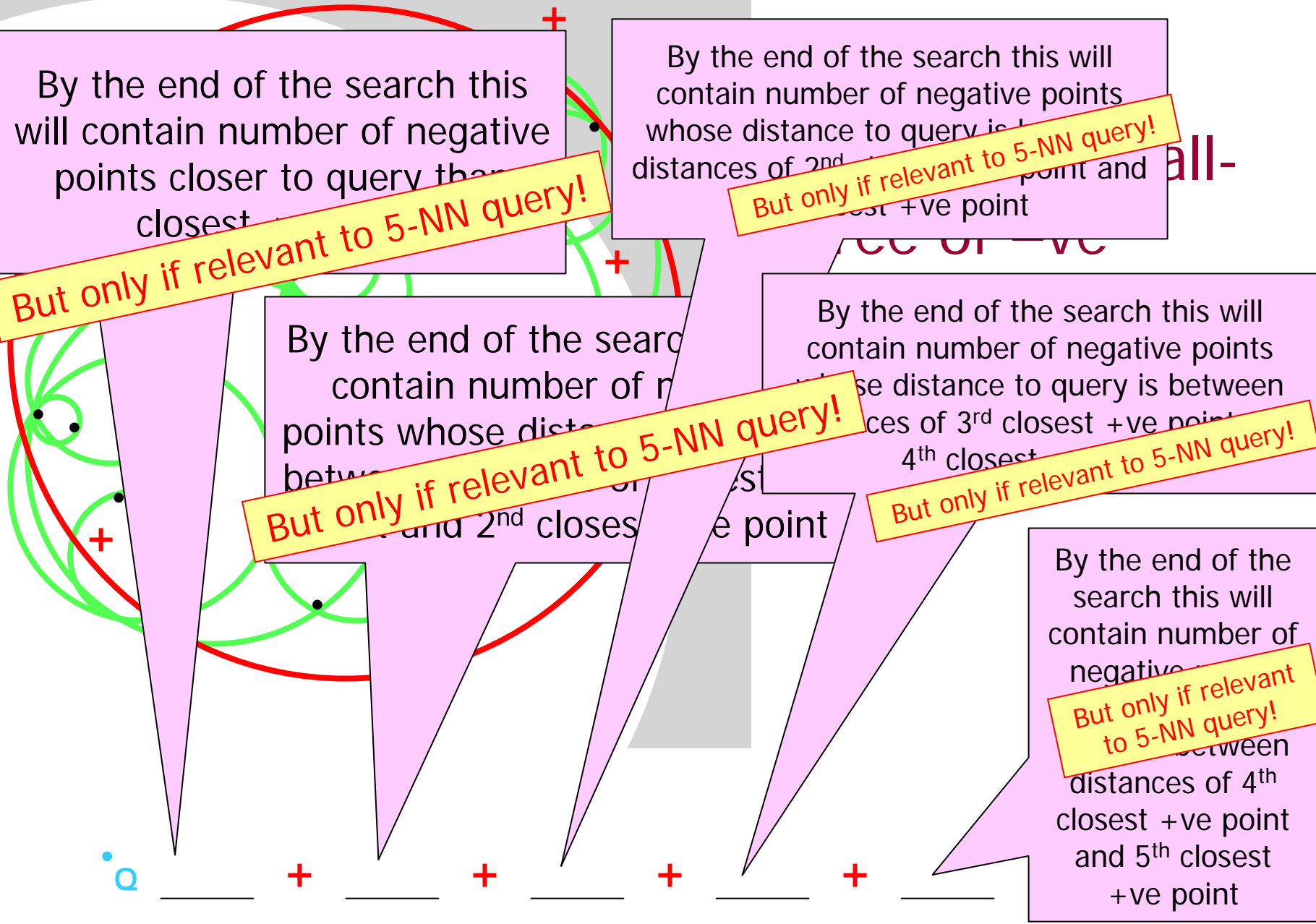
But only if relevant to 5-NN query!

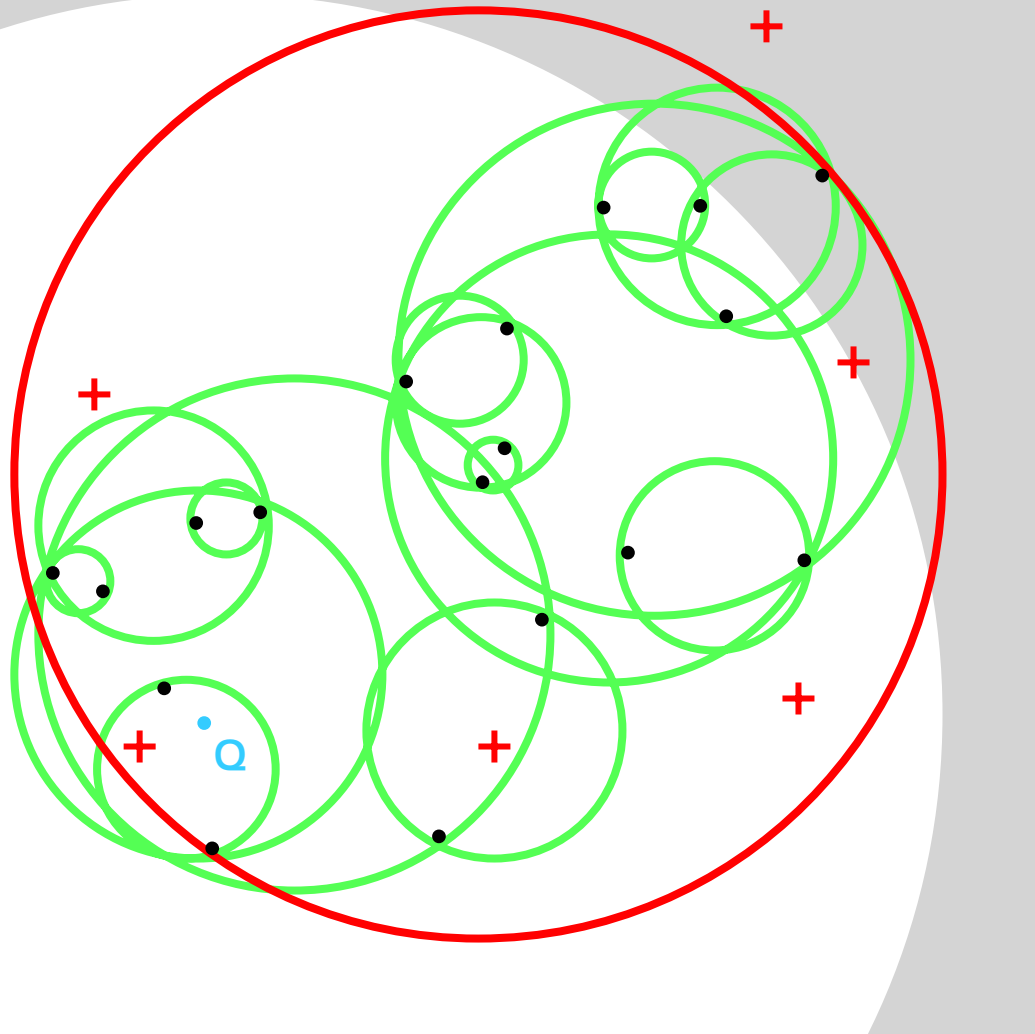
By the end of the search this will contain number of negative points whose distance to query is between distances of 4th and 5th closest +ve point

But only if relevant to 5-NN query!

By the end of the search this will contain number of negative points whose distance to query is between distances of 5th closest +ve point and 6th closest +ve point

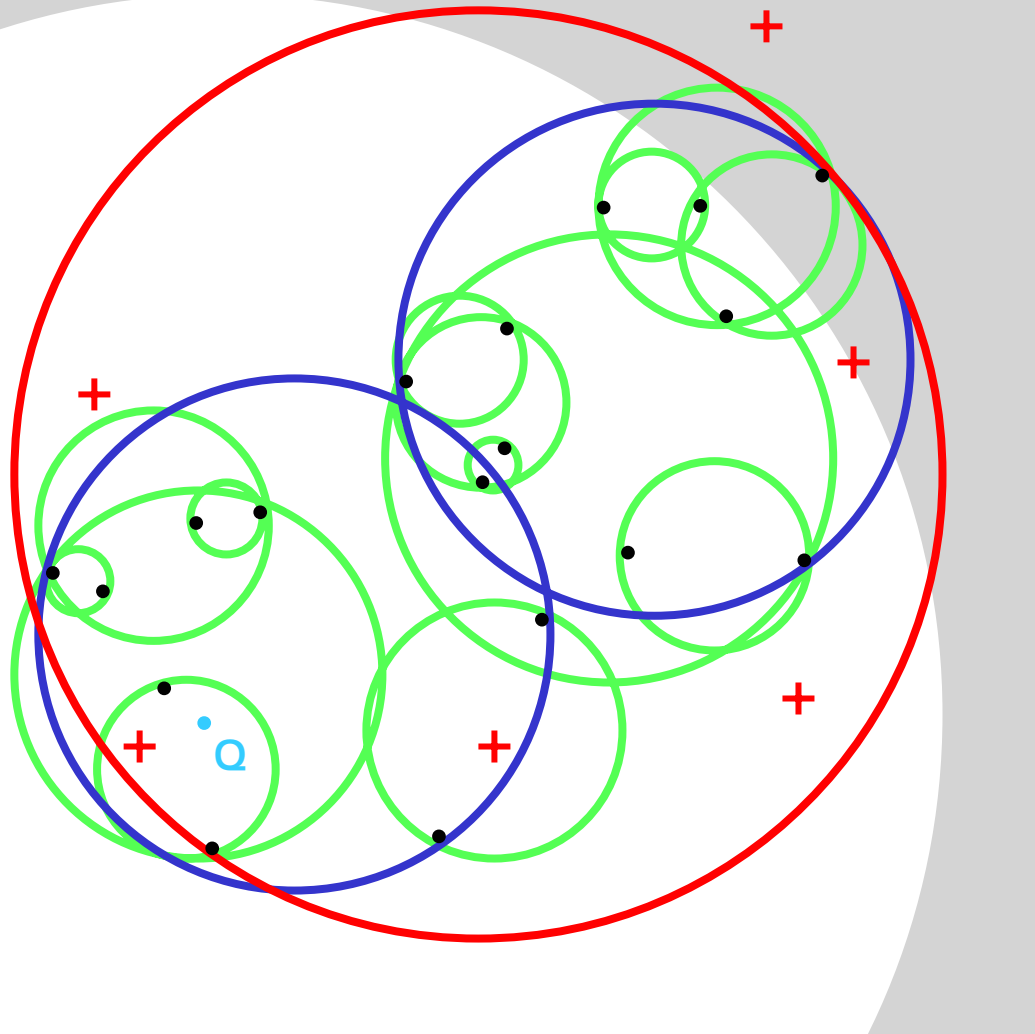
But only if relevant to 5-NN query!



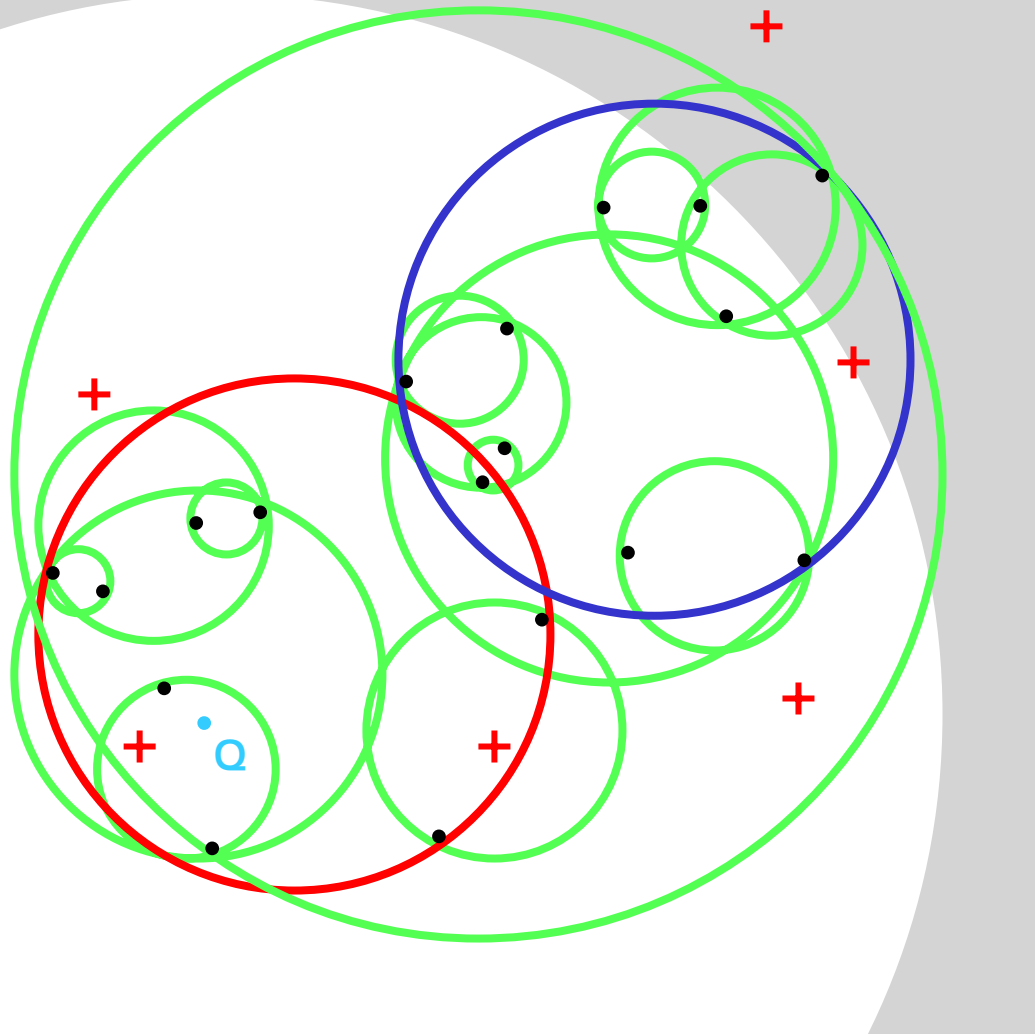


Search the ball-tree of -ve points starting at the root.

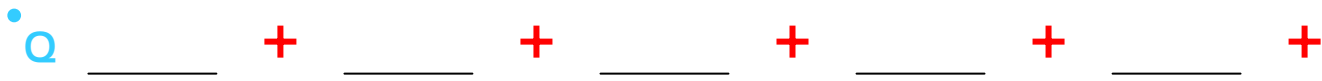
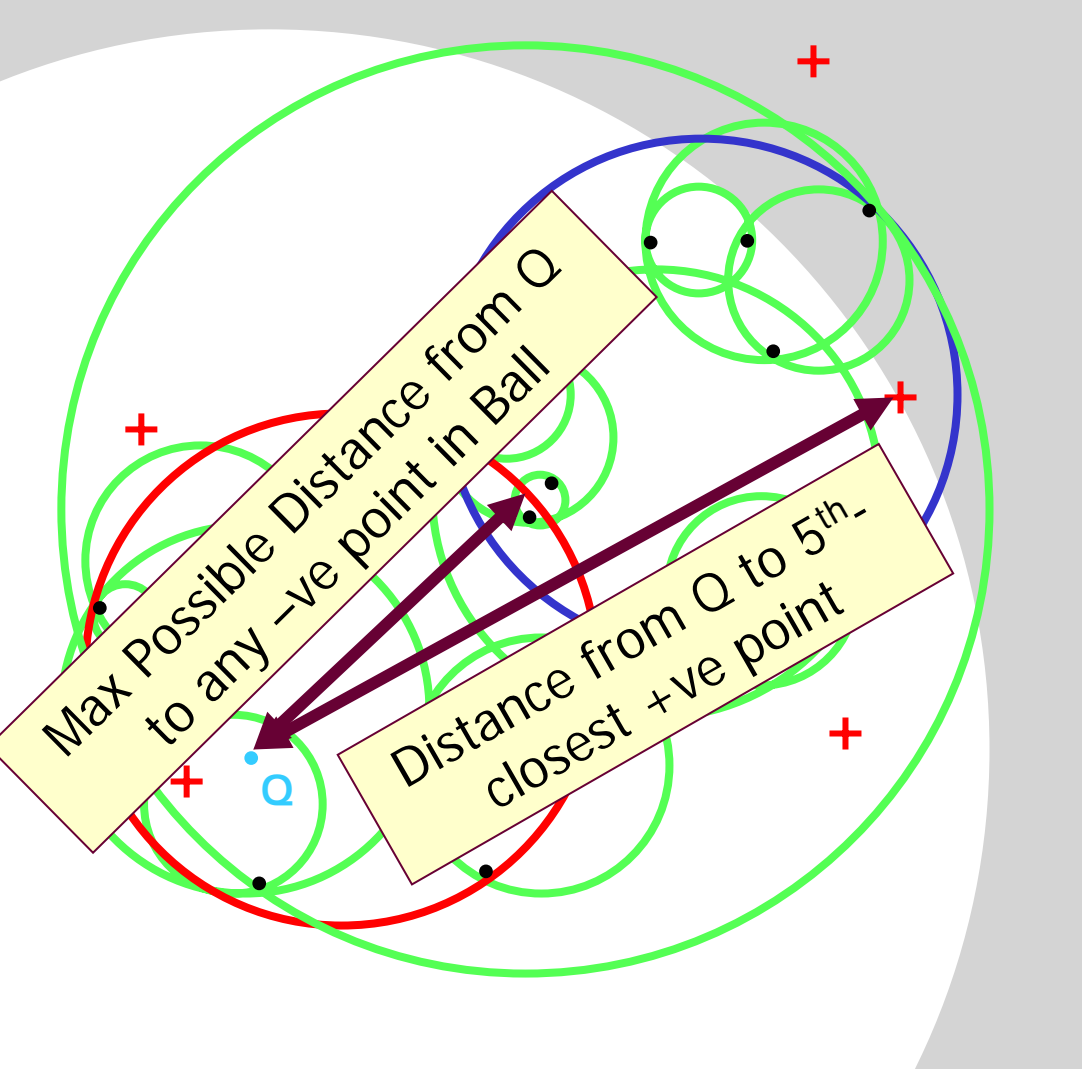


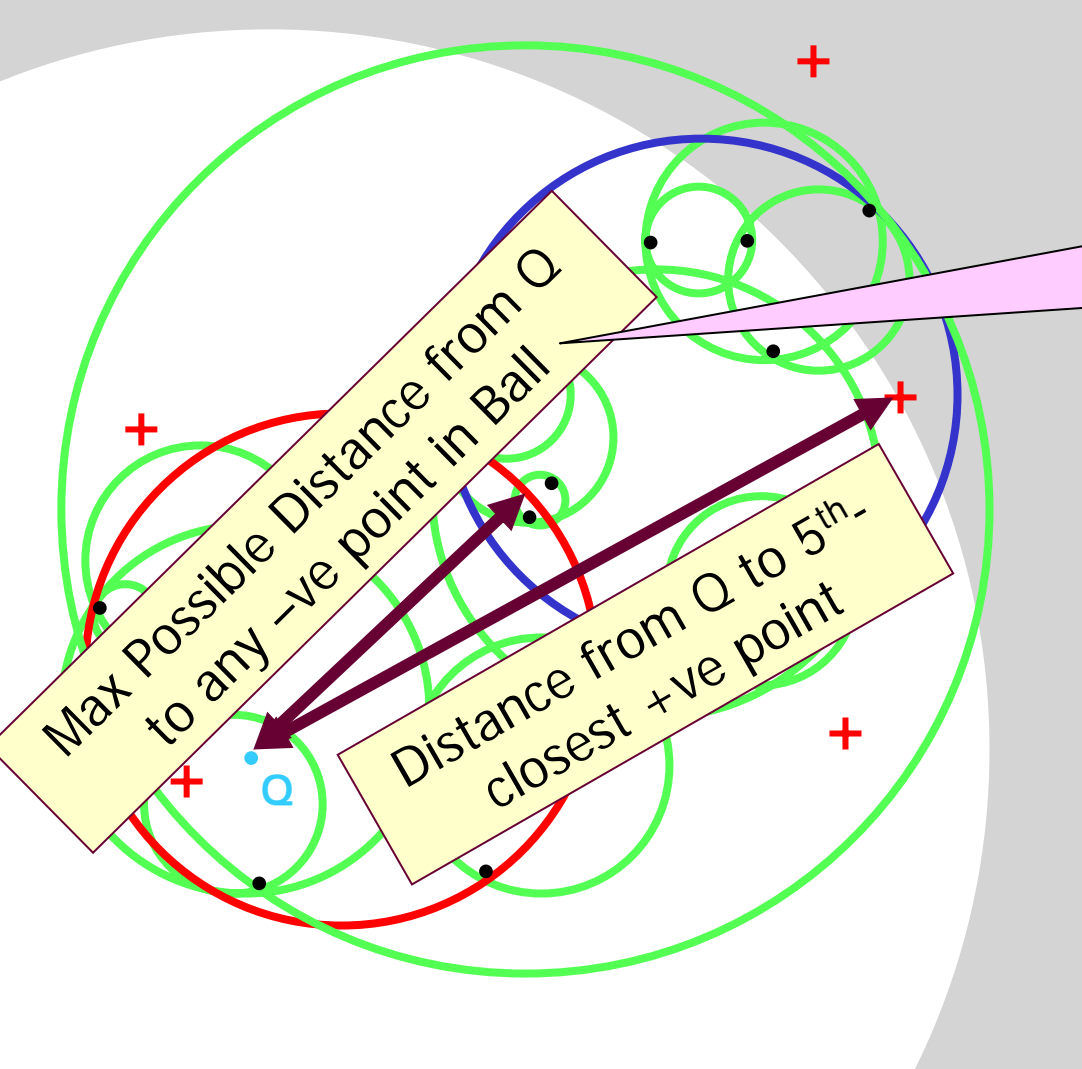


\circ q $+$ $+$ $+$ $+$ $+$ $+$

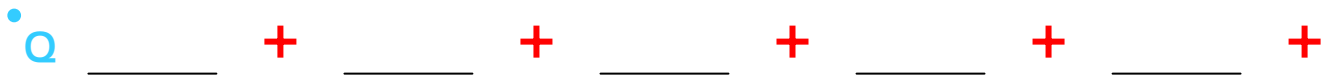


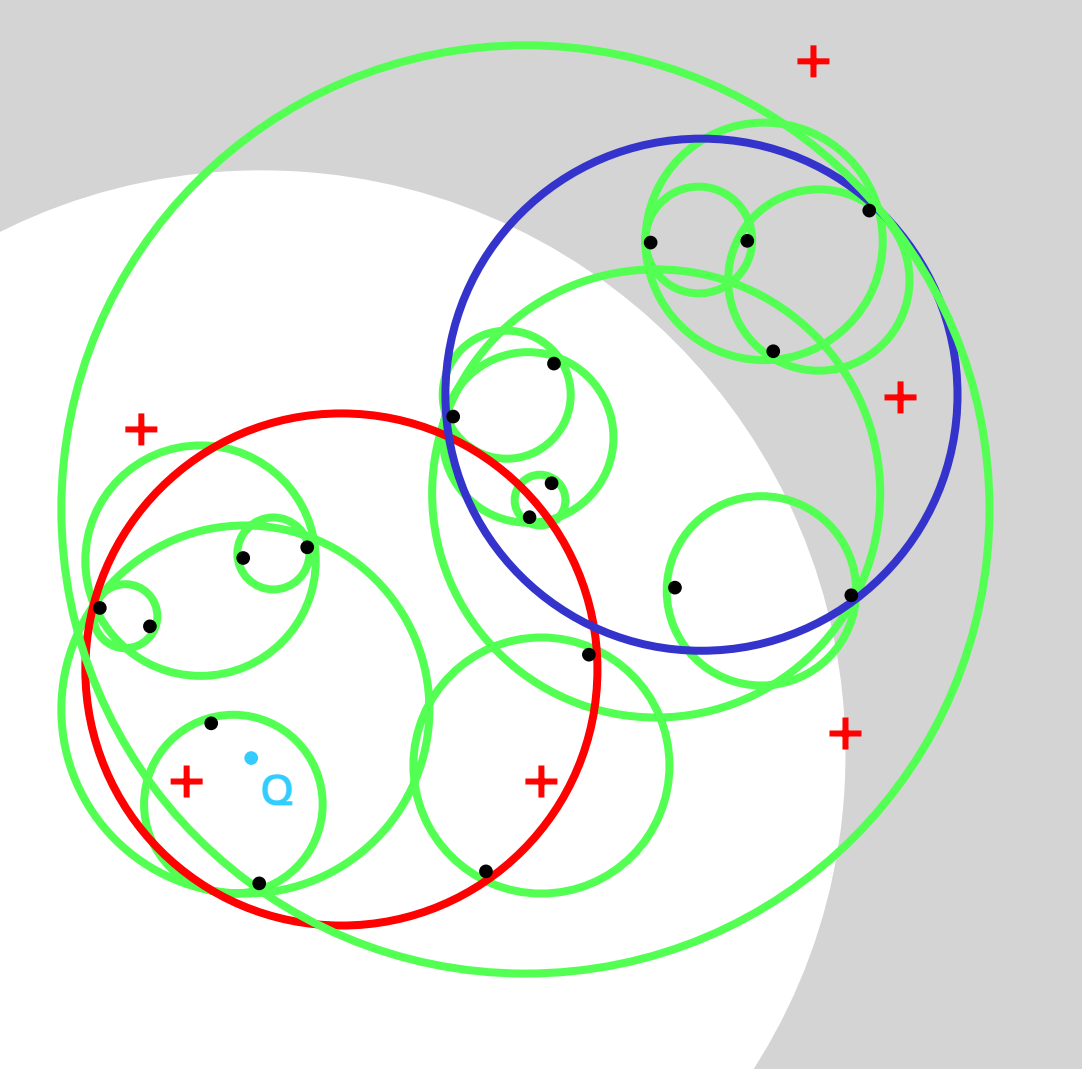
\circ Q $+$ $+$ $+$ $+$ $+$ $+$



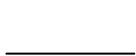


...and there are eight points in the ball





Q



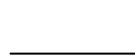
+



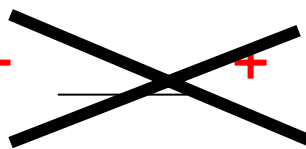
+



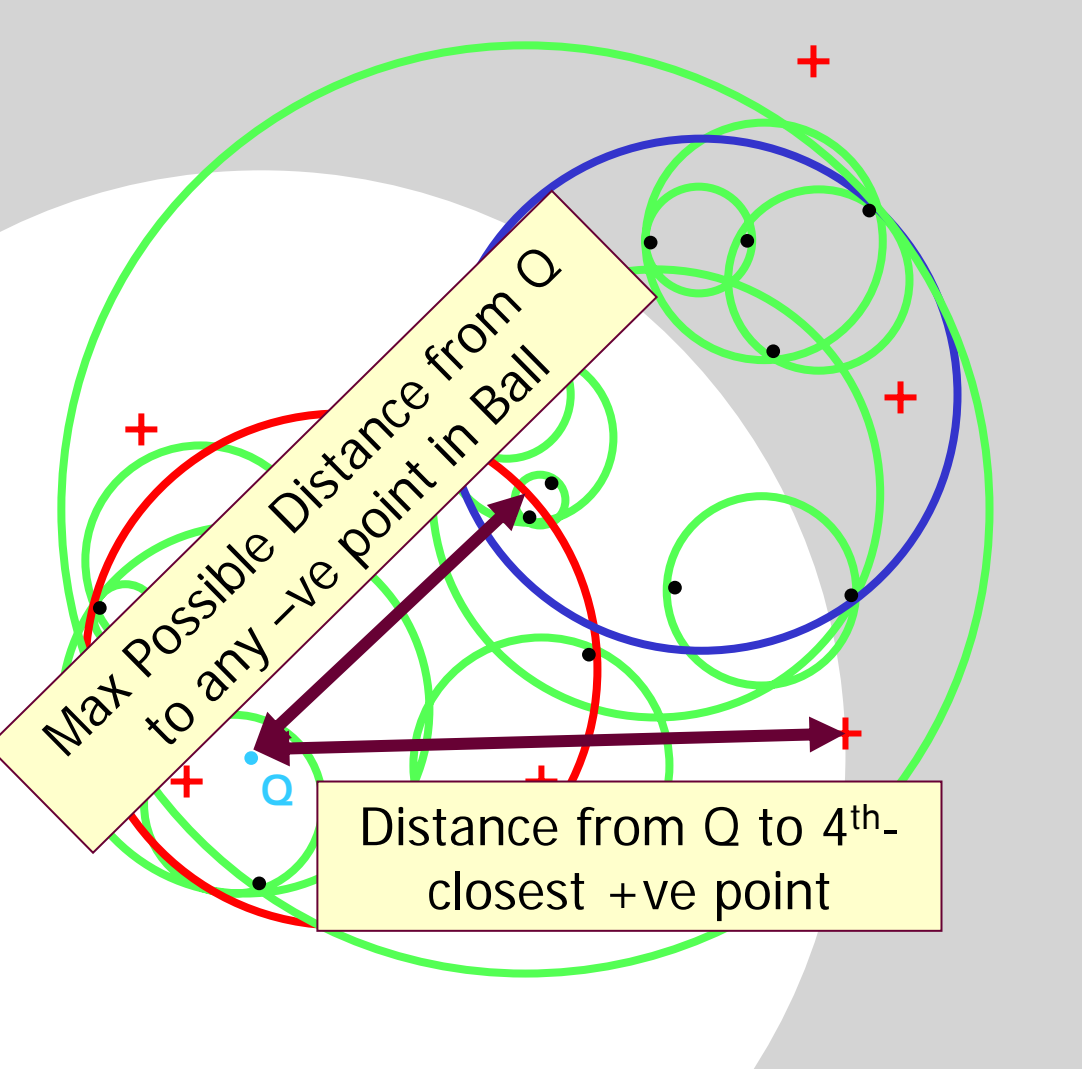
+

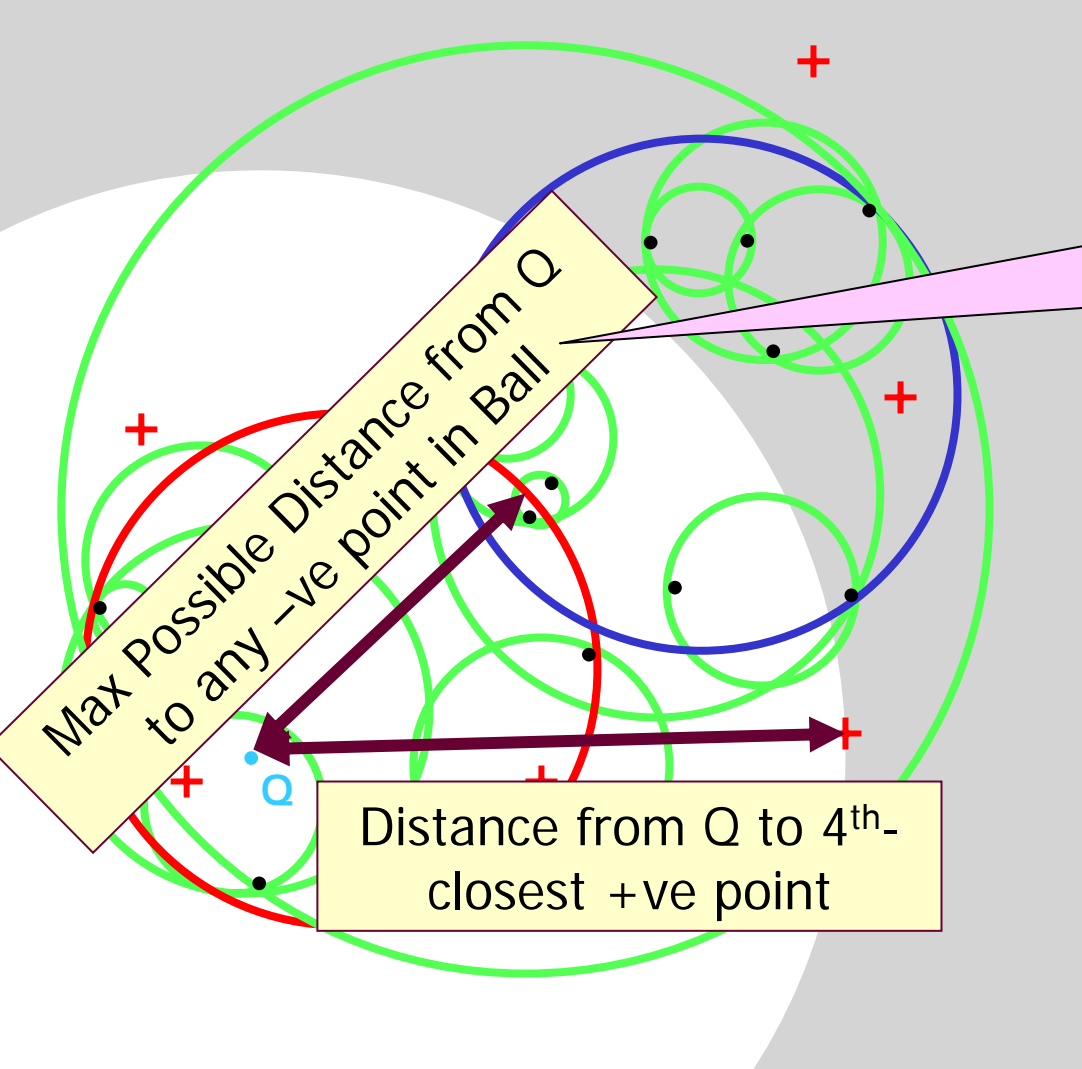


+



+



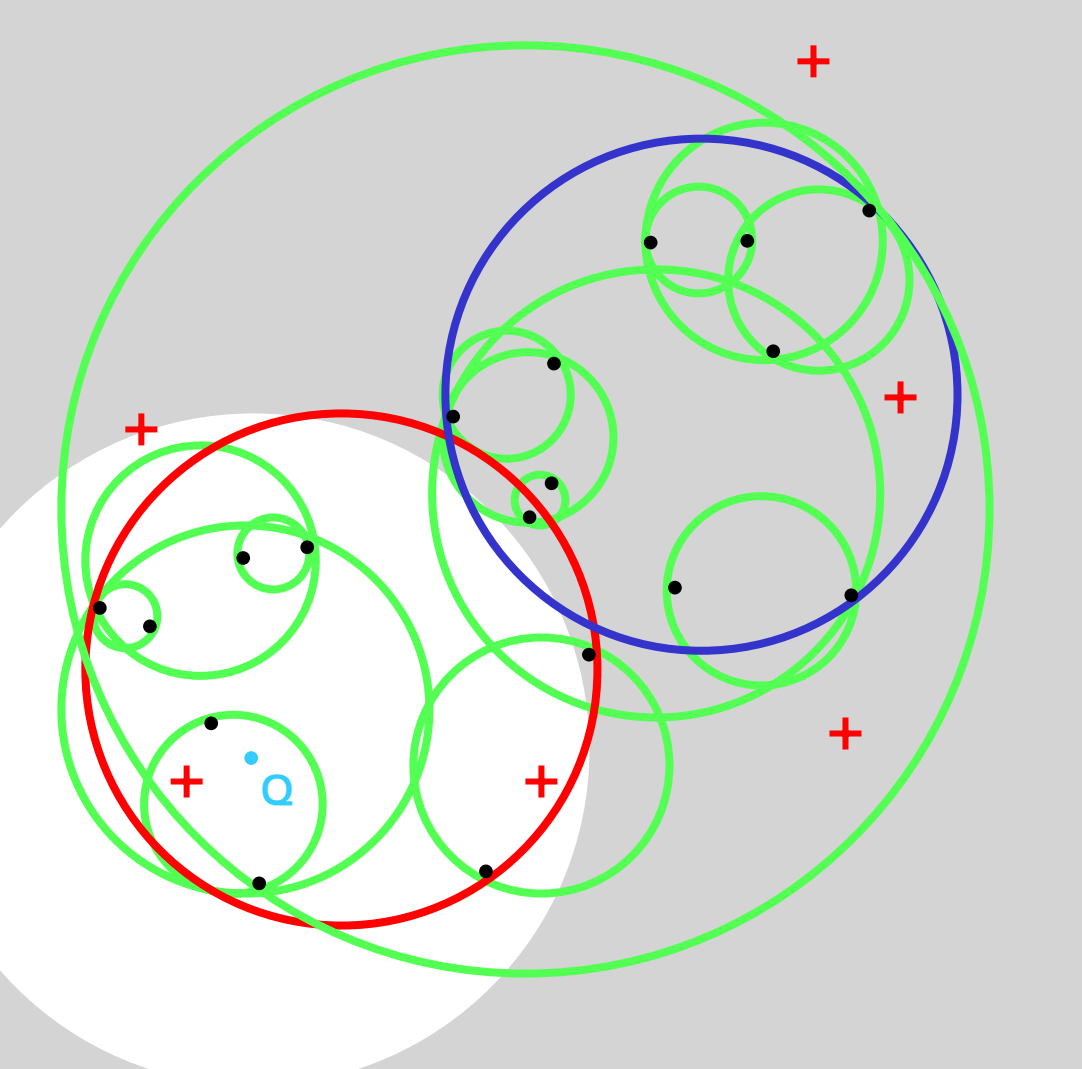


...and there are eight points in the ball

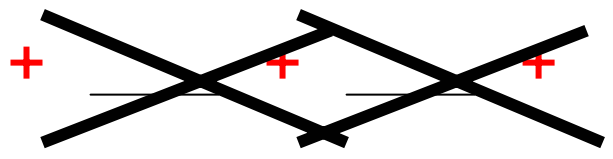
Distance from Q to 4th-closest +ve point

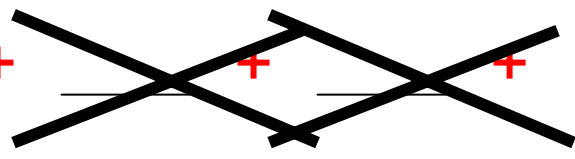
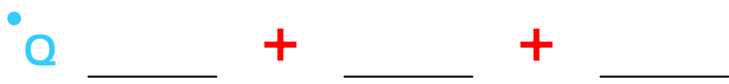
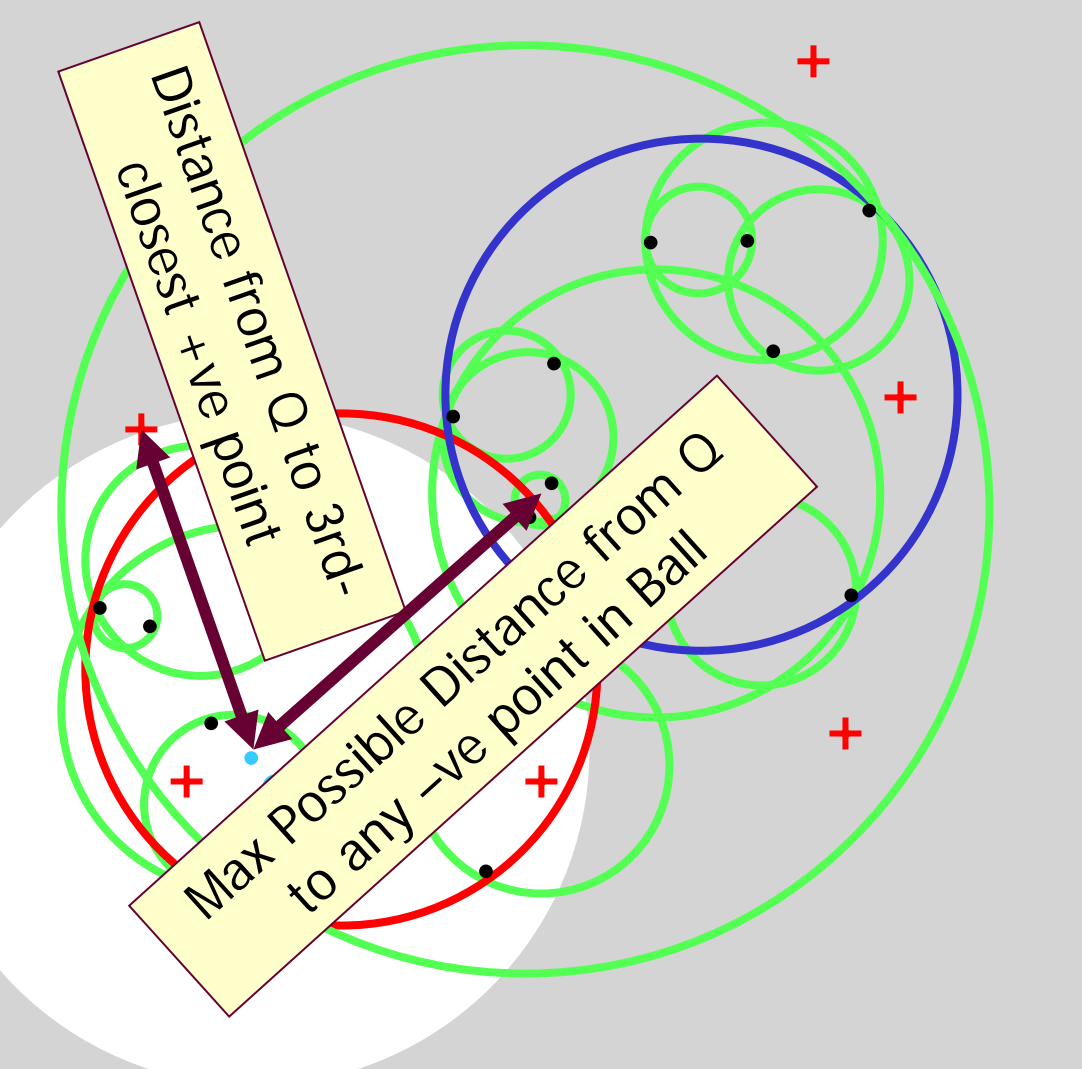
Max Possible Distance from Q to any -ve point in Ball

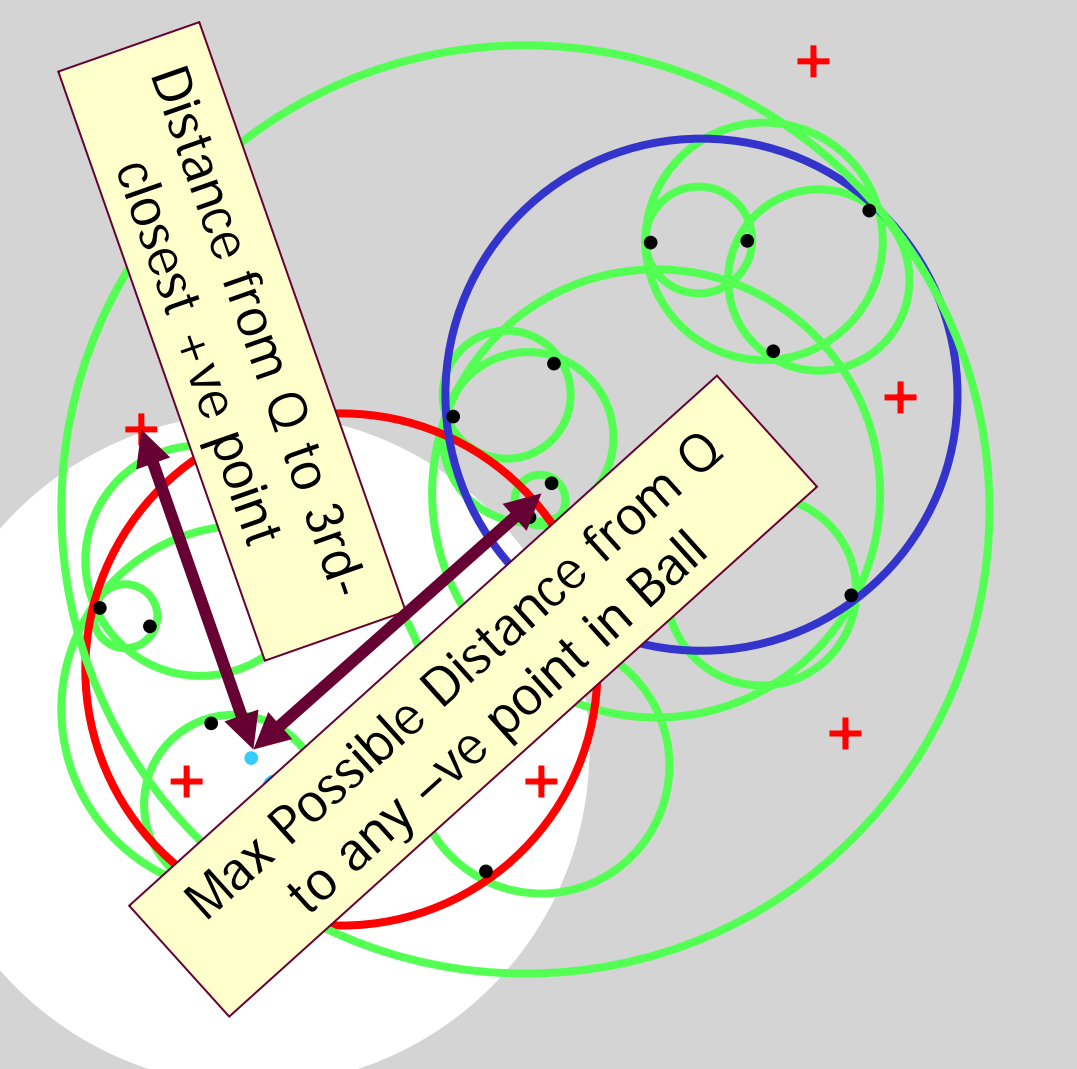




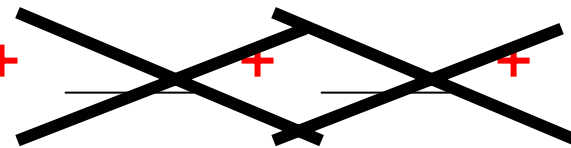
Q _____ + _____ + _____

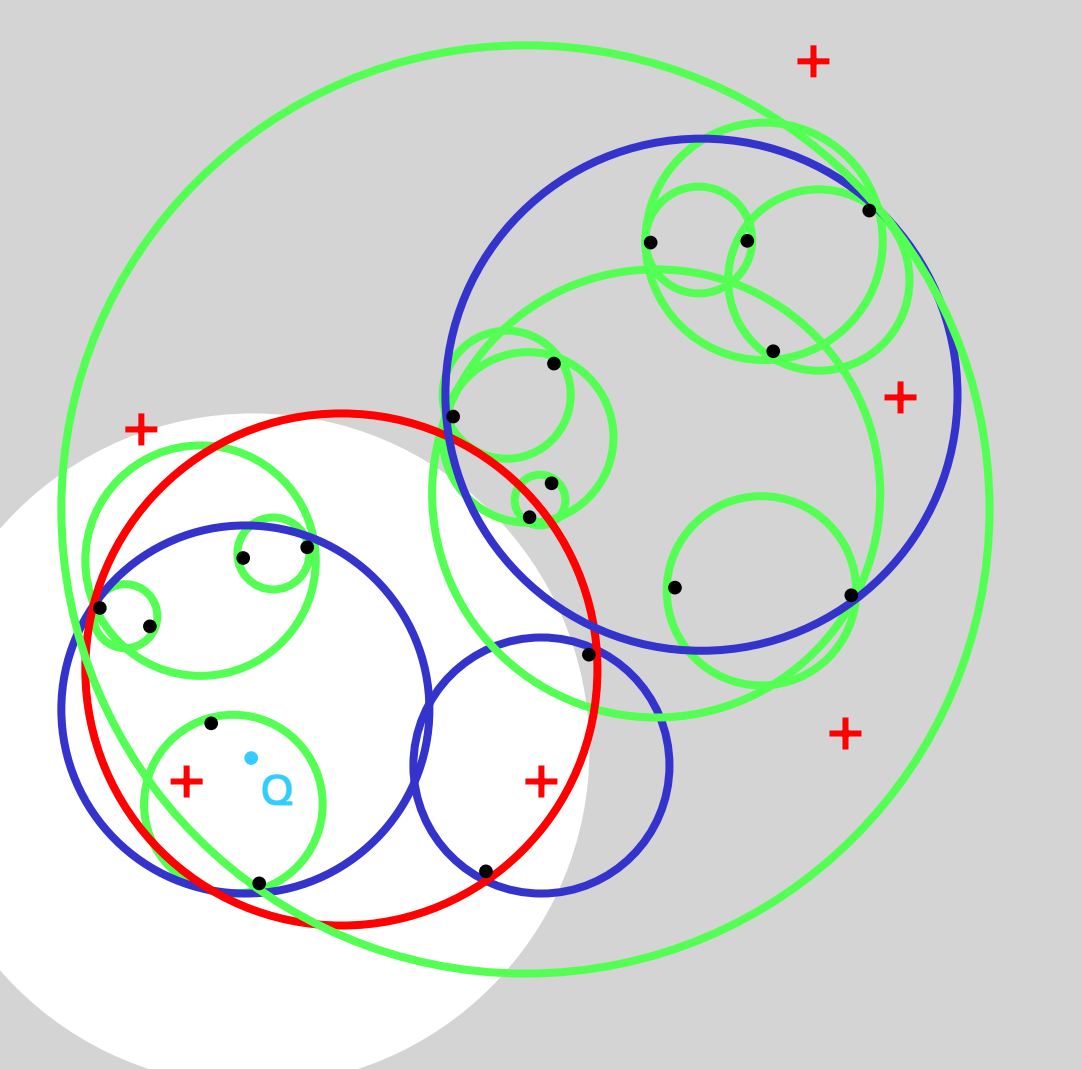




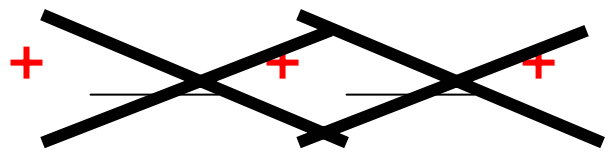


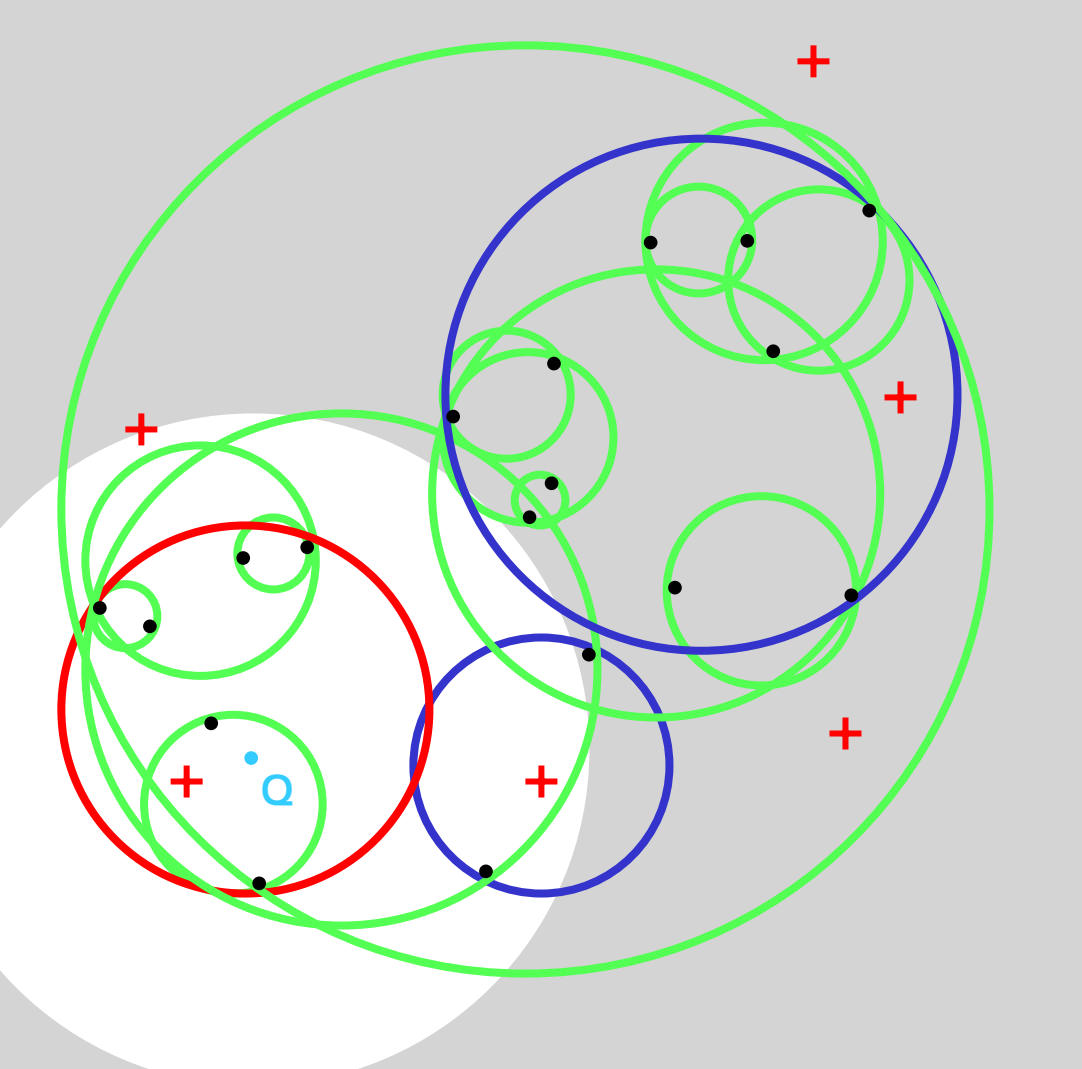
No prune!



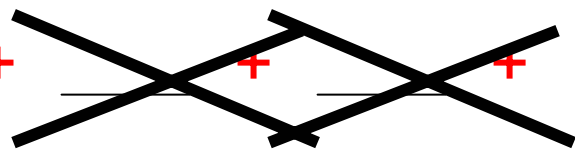


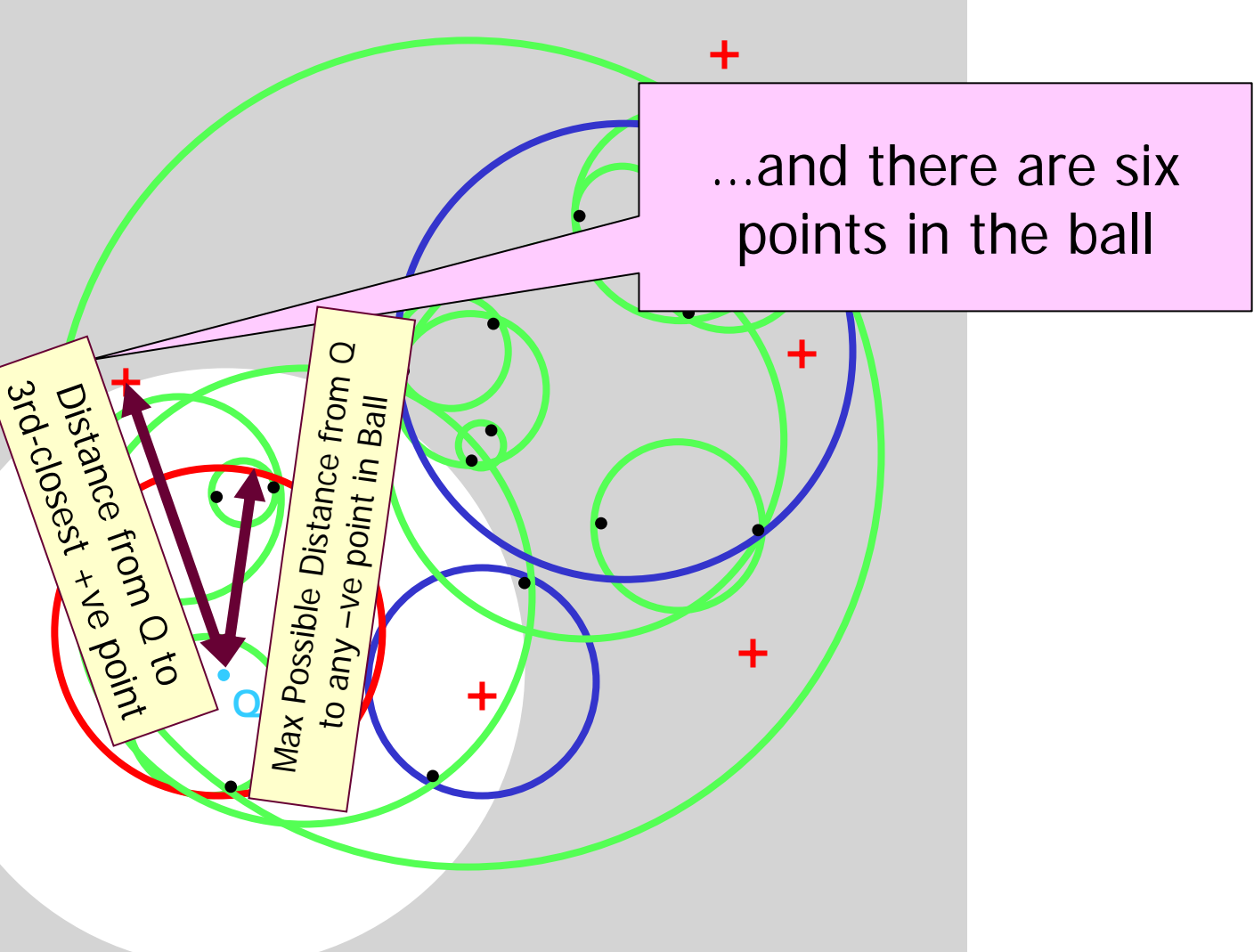
Q _____ + _____ + _____



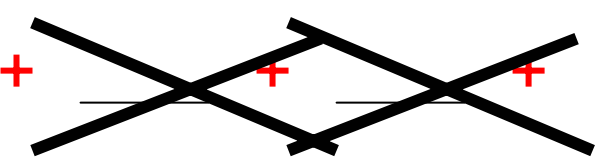
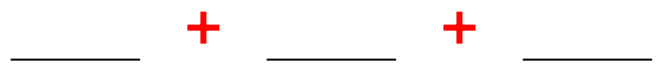


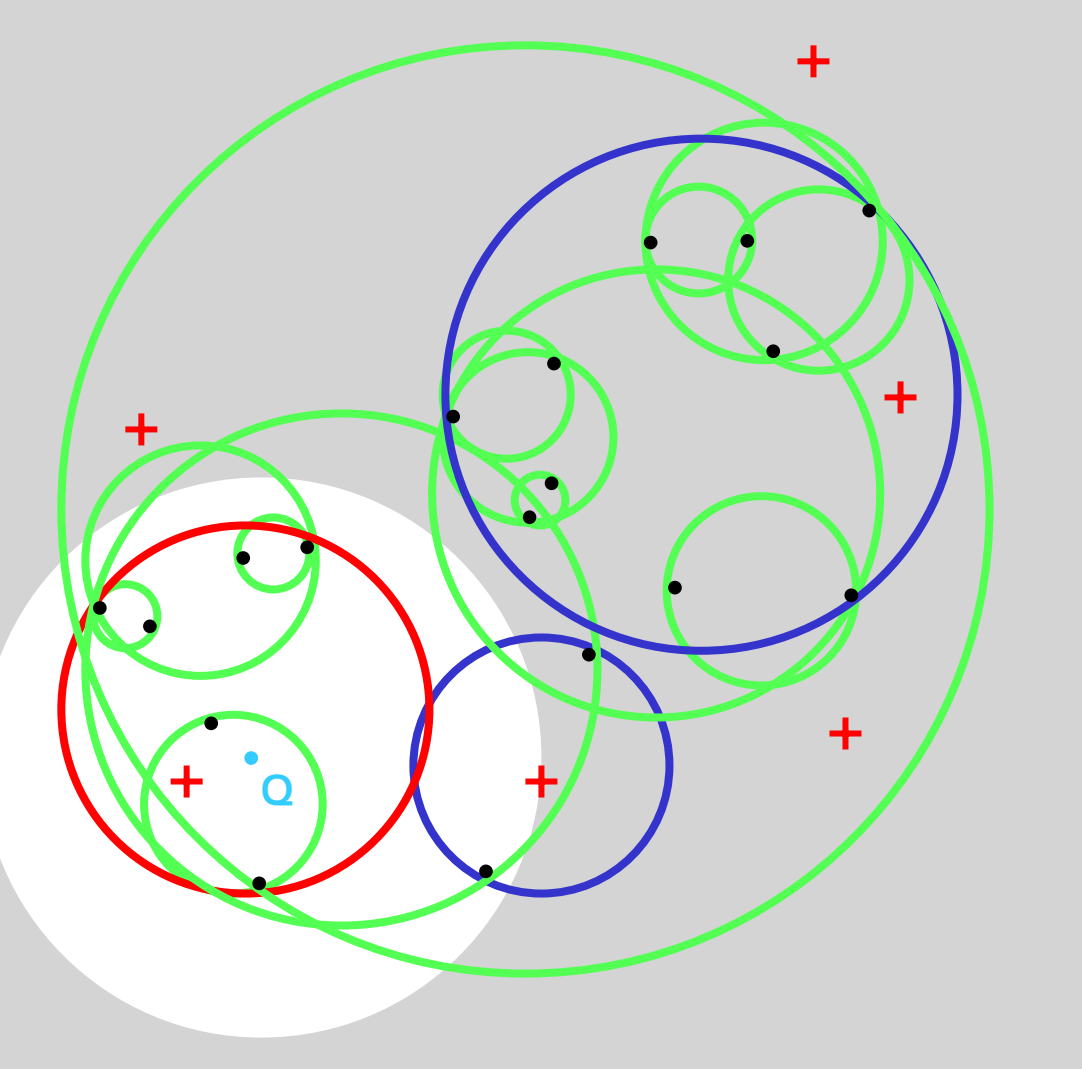
Q _____ + _____ + _____





Q

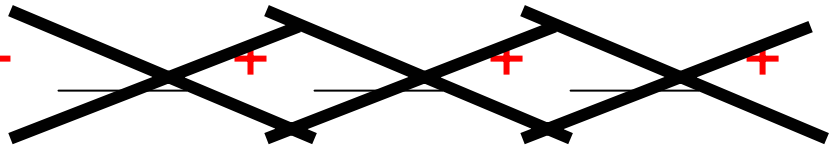




Q

+

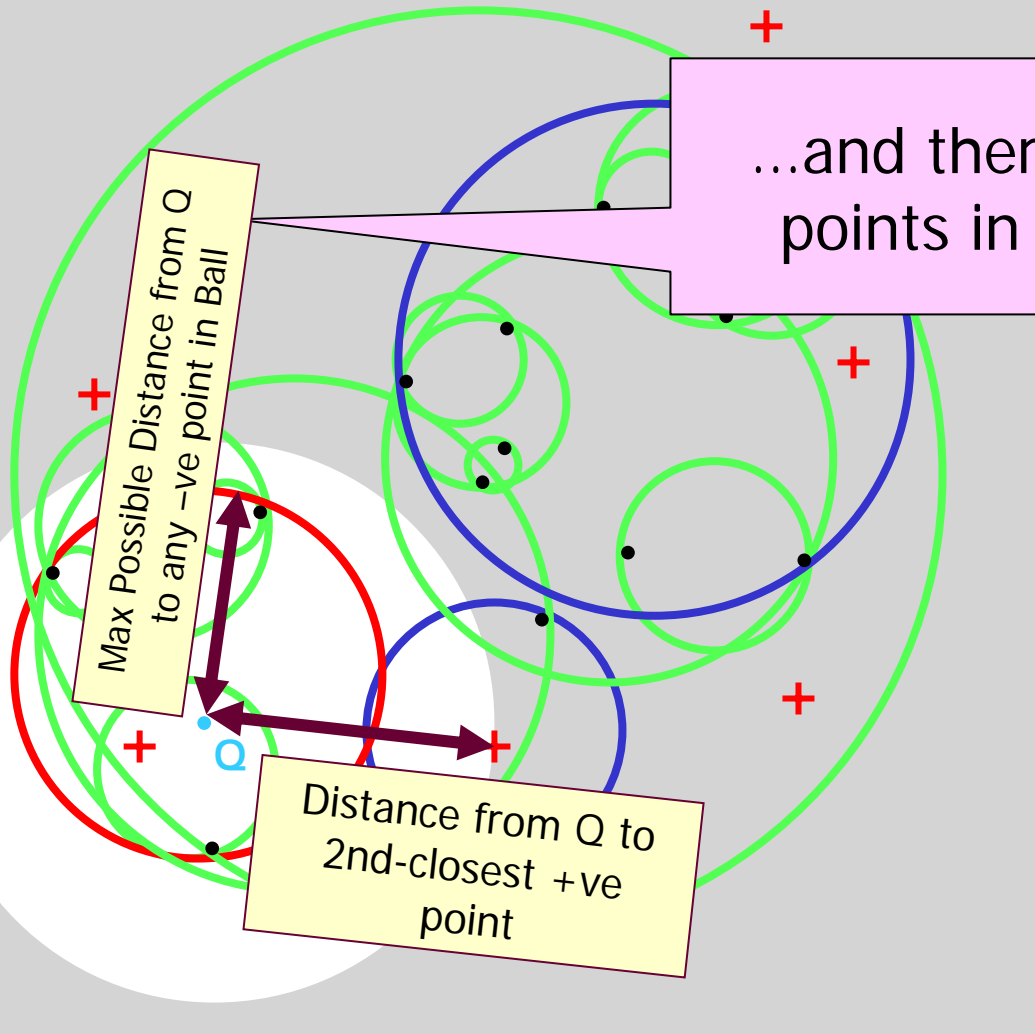
+



+

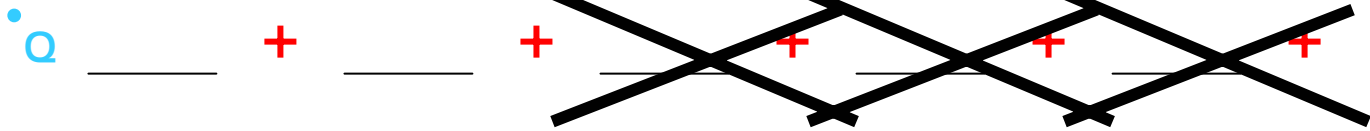
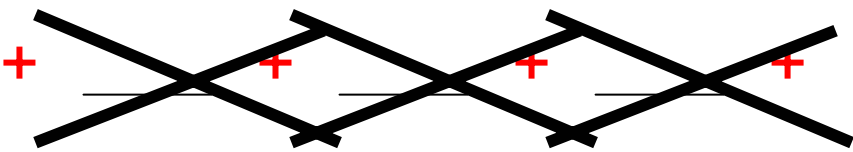
+

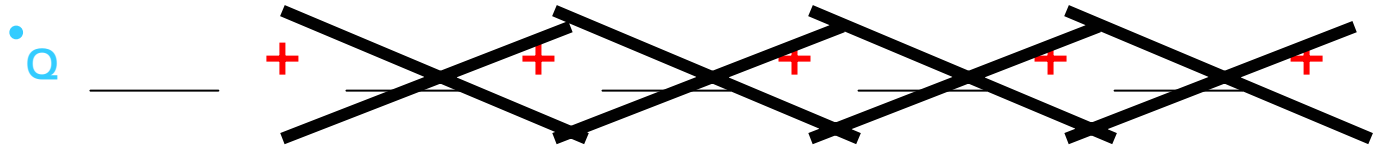
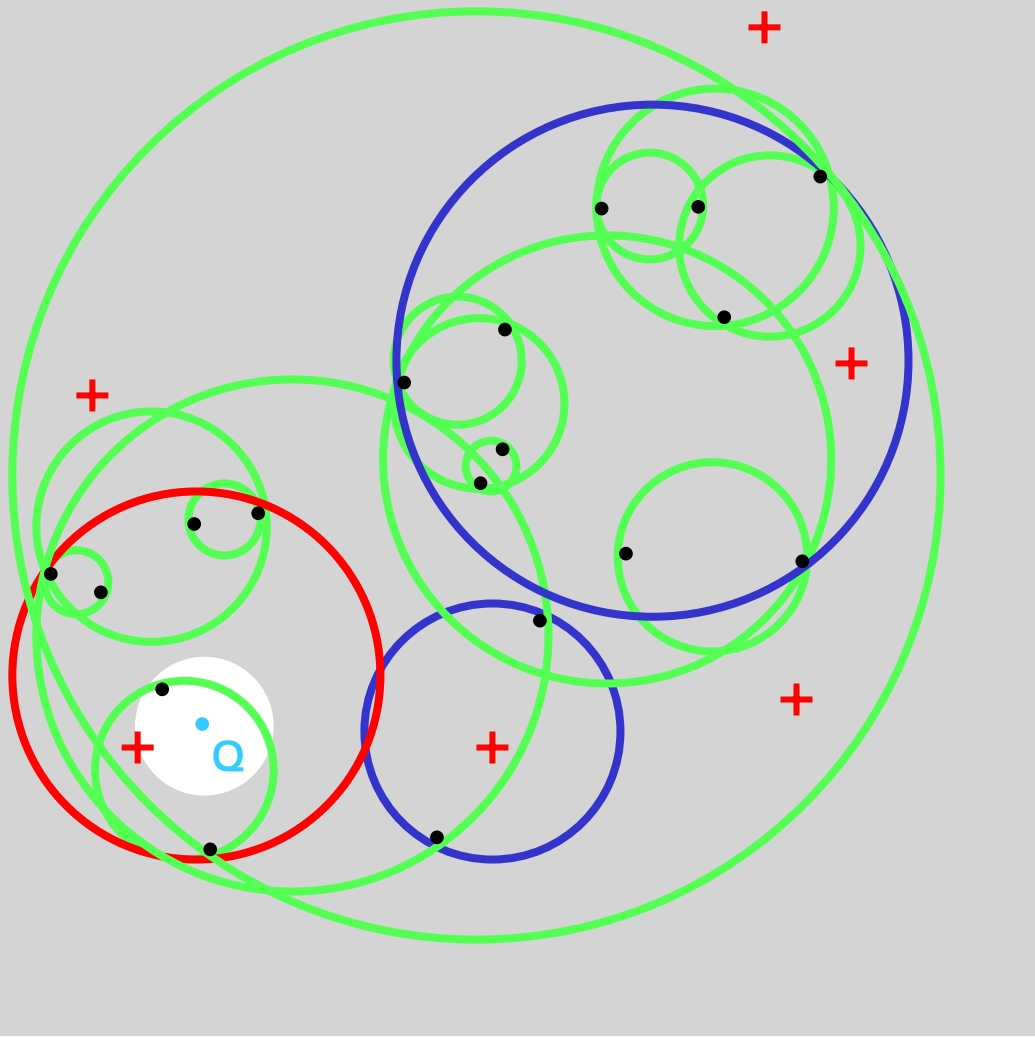
+

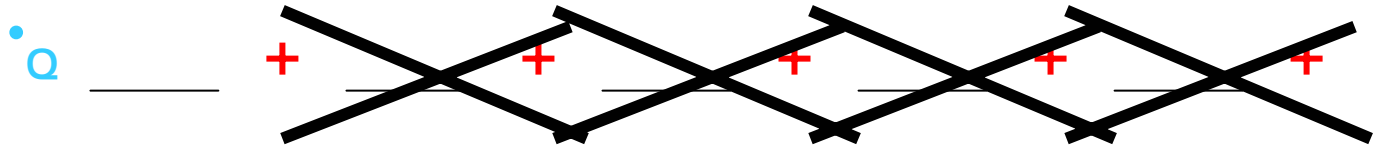
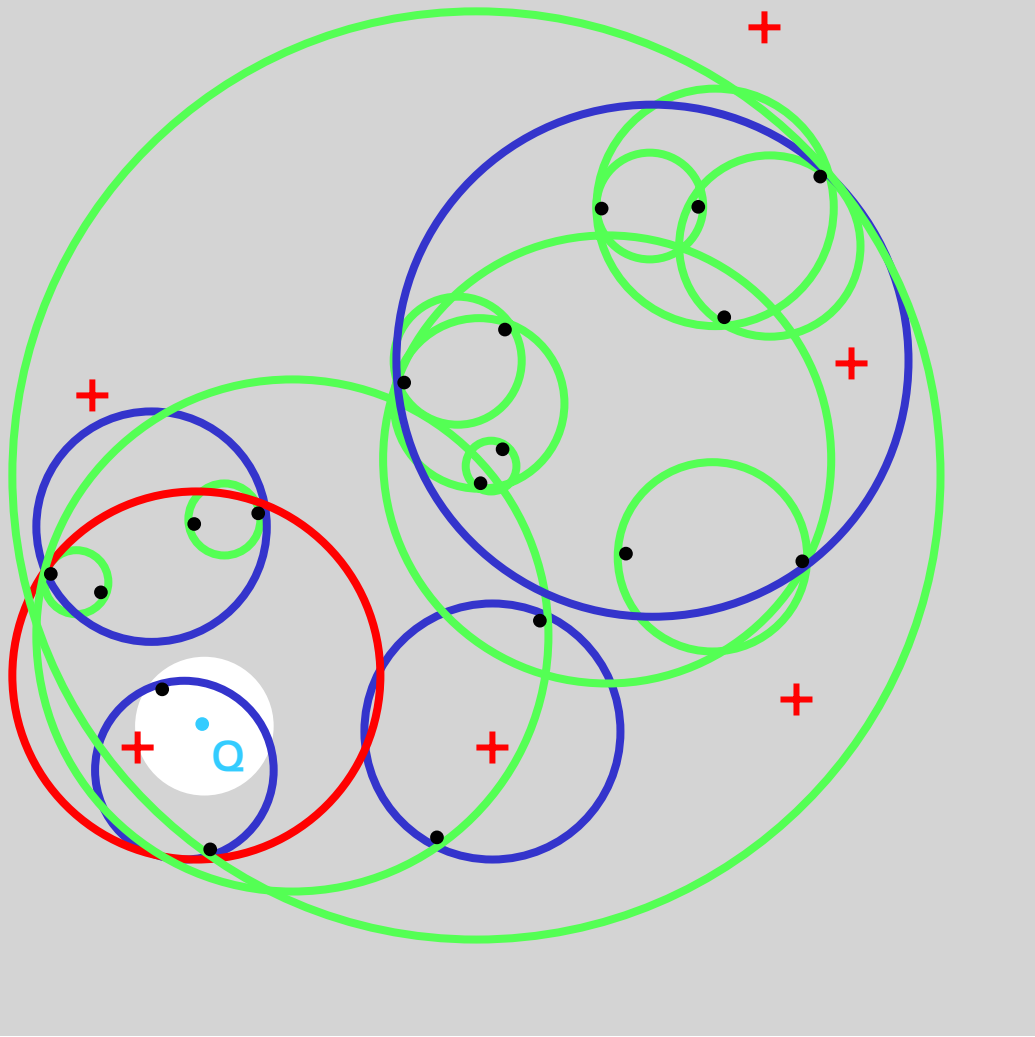


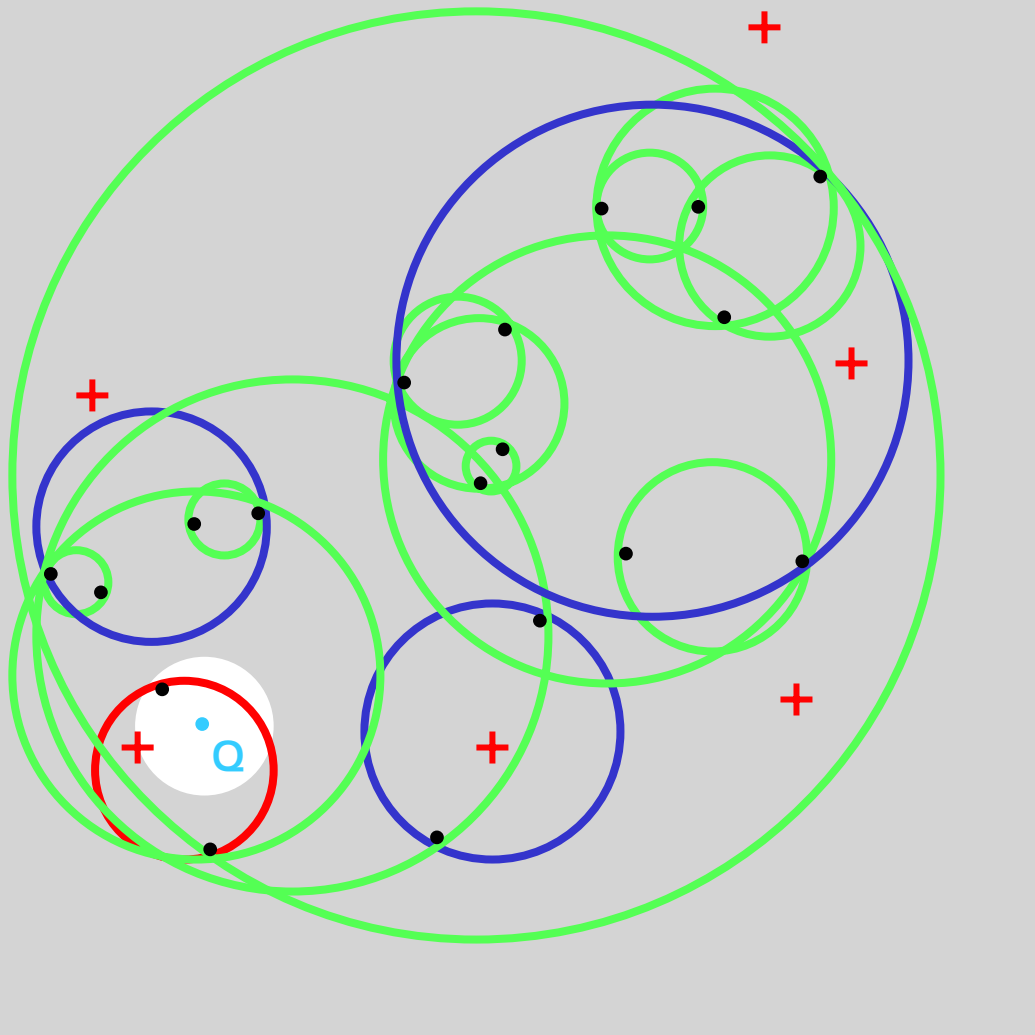
...and there are six points in the ball

Q

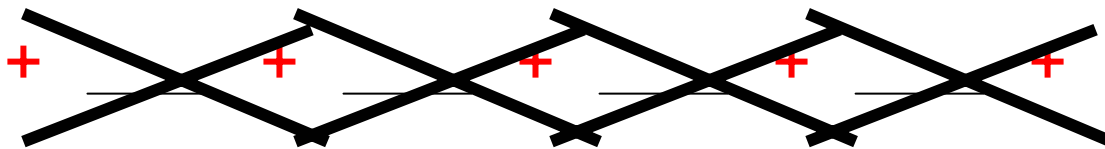


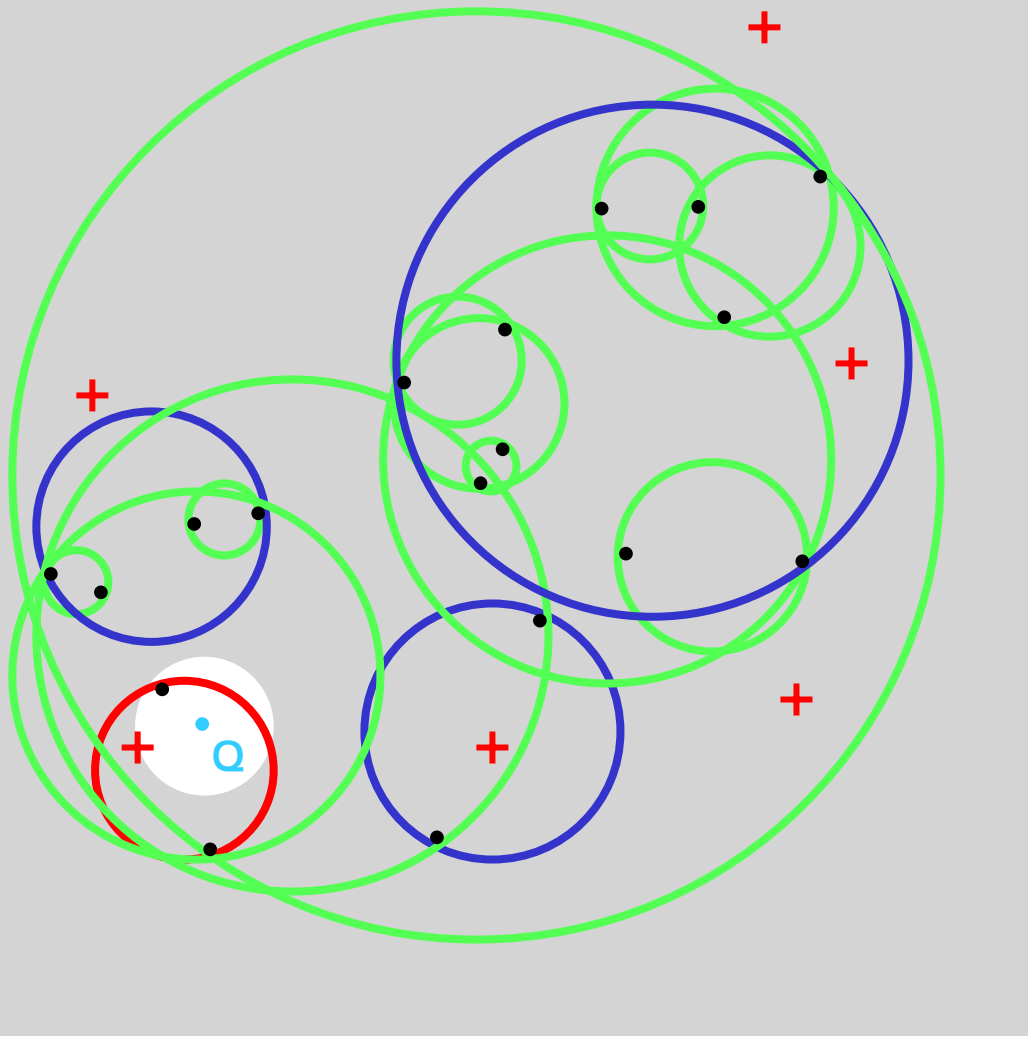




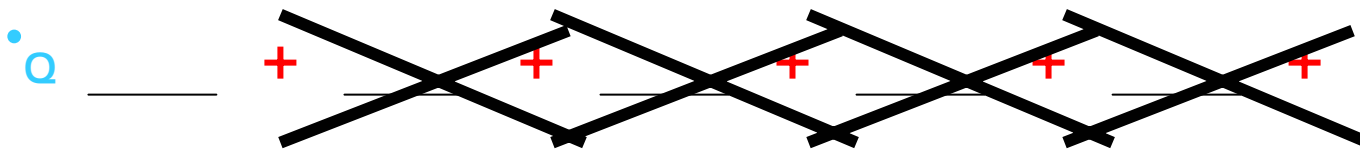


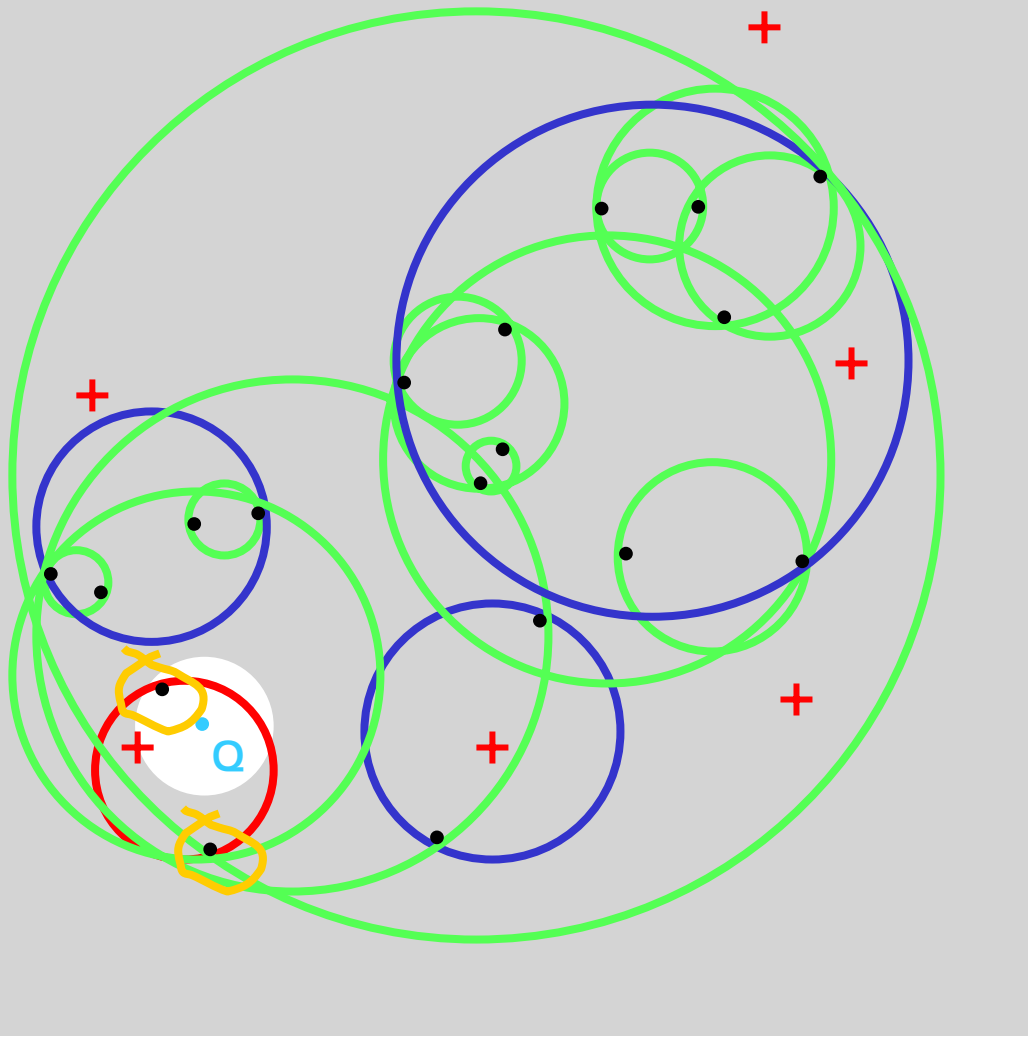
Q





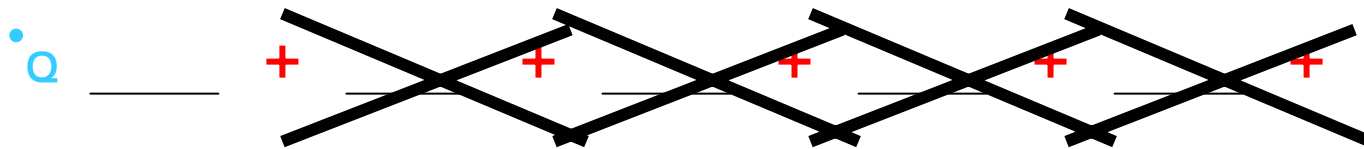
No prune.

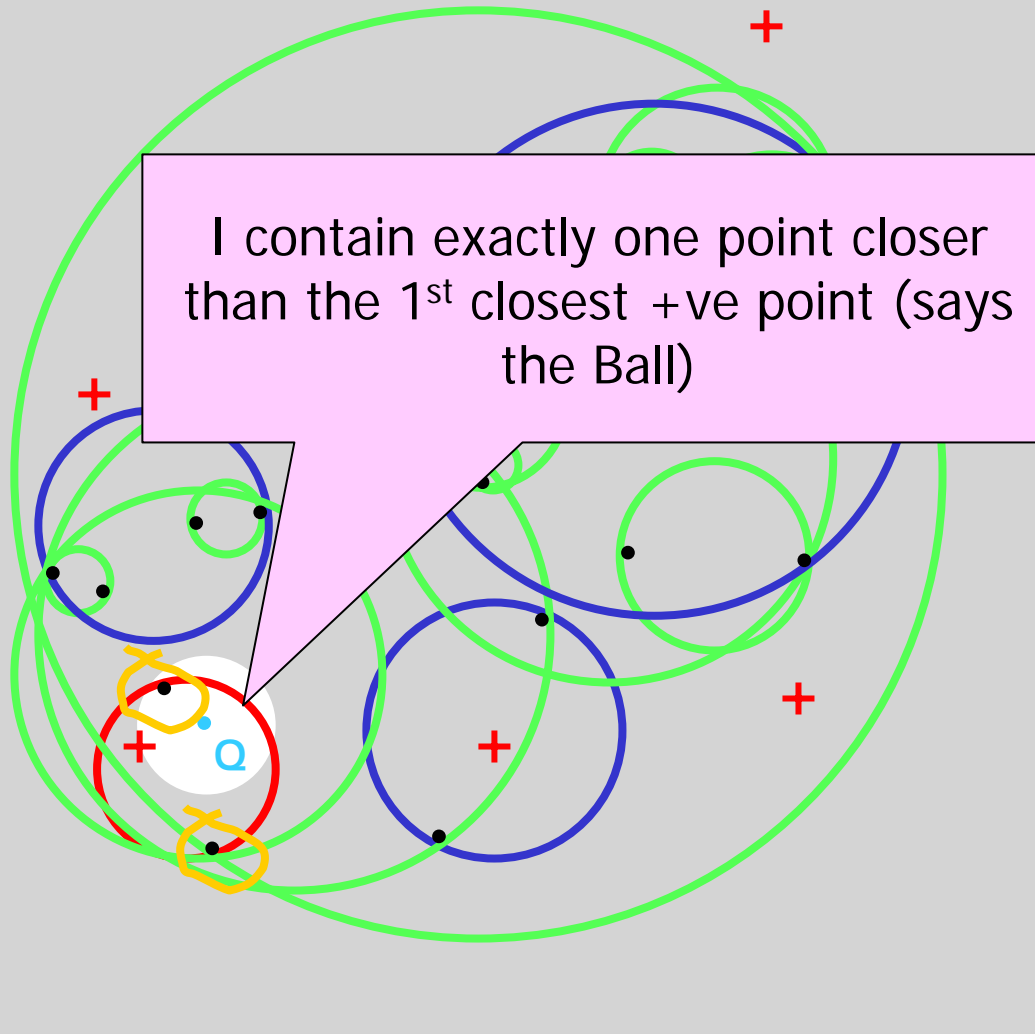




No prune.

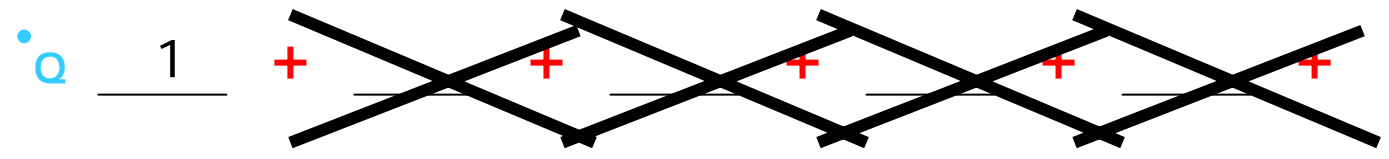
Ball is leaf
so explore
its points

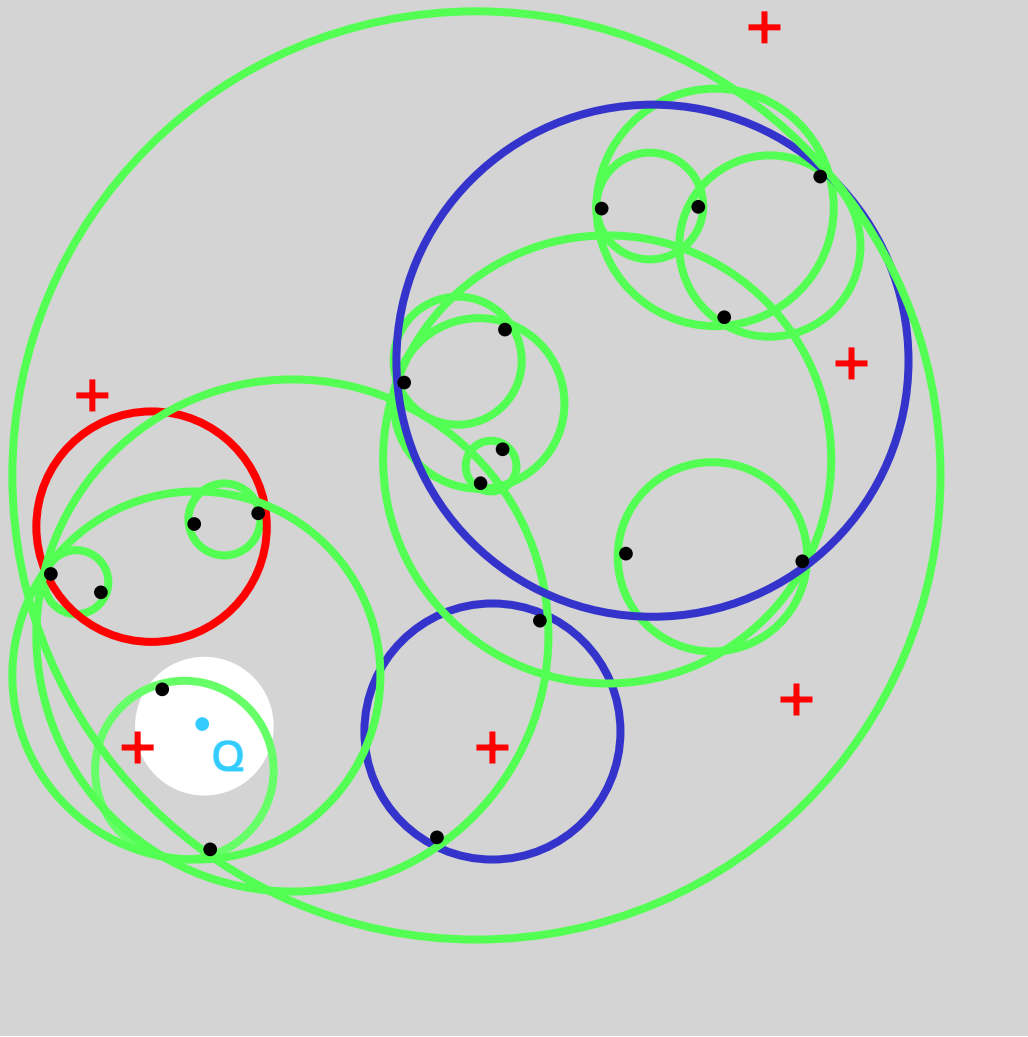




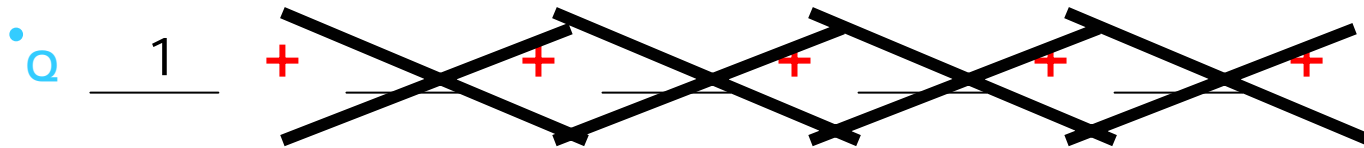
No prune.

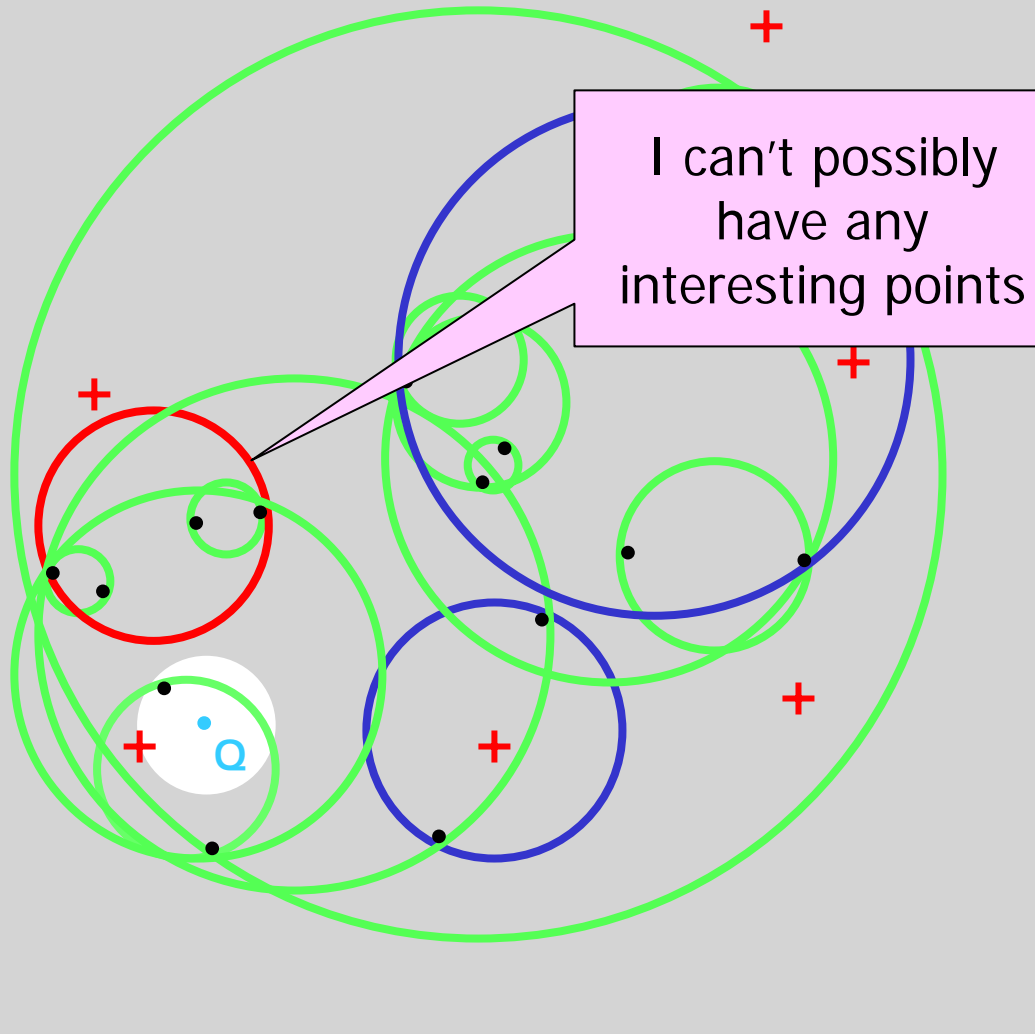
Ball is leaf
so explore
its points





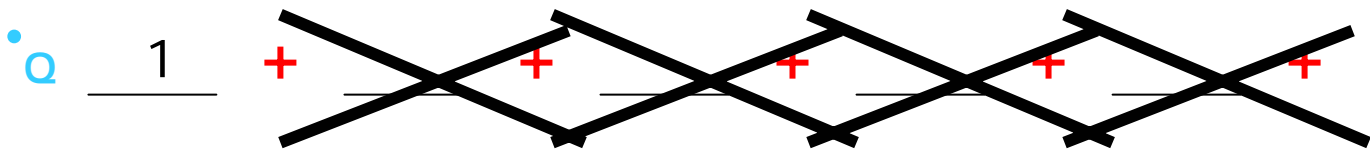
Return and
try other
sibling

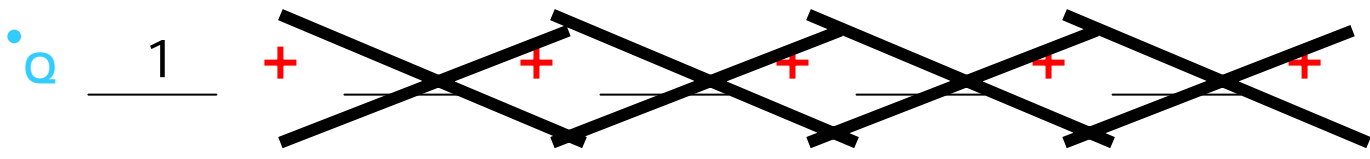
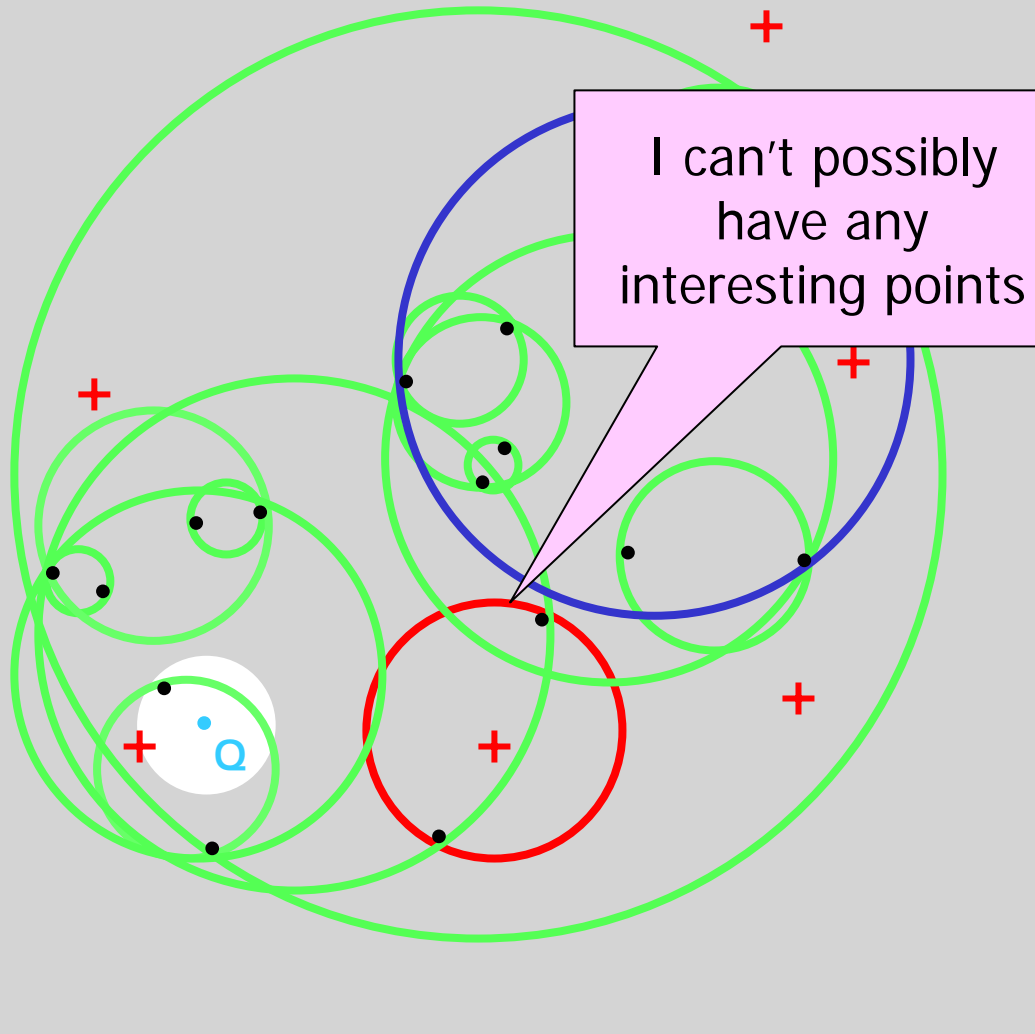


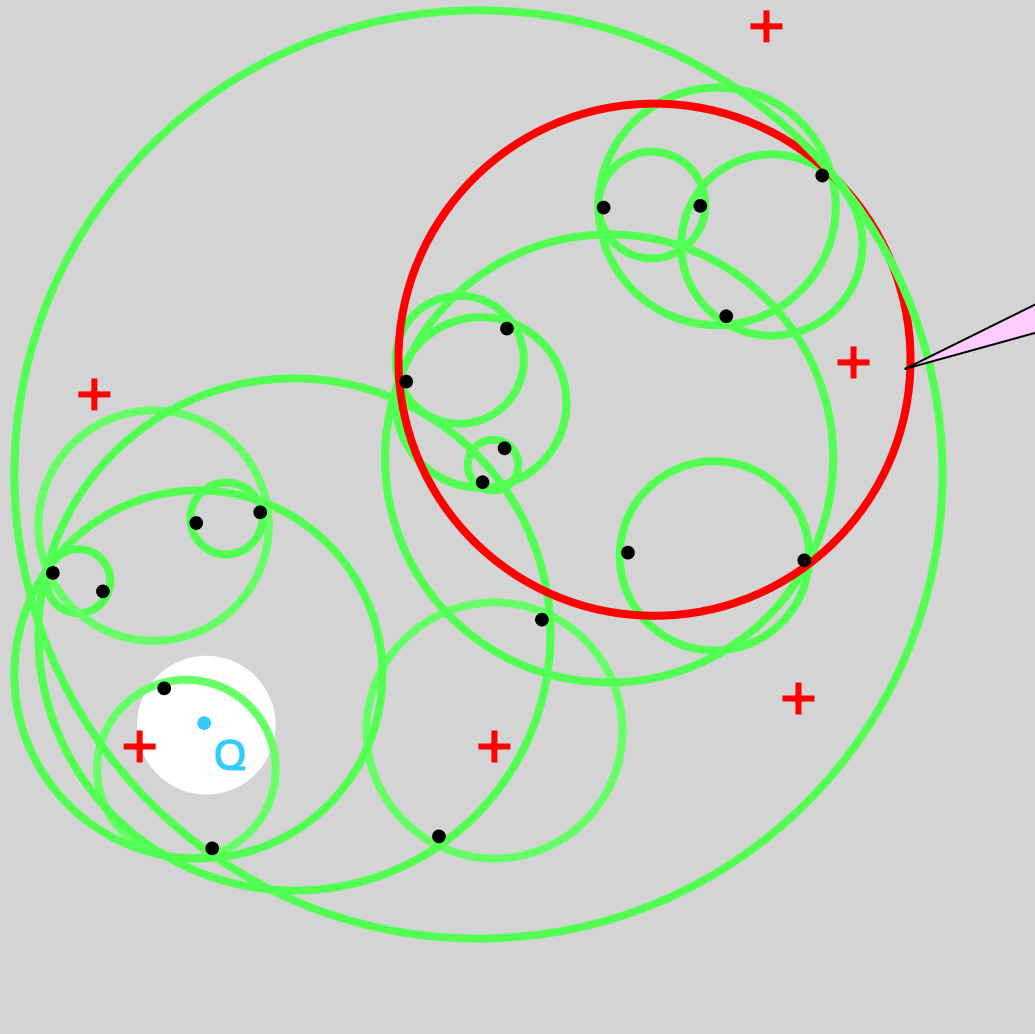


I can't possibly have any interesting points

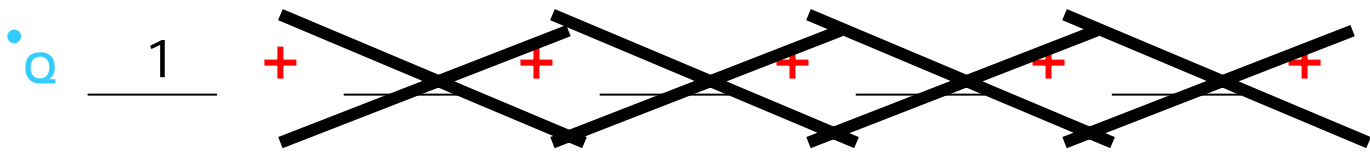
Return and try other sibling



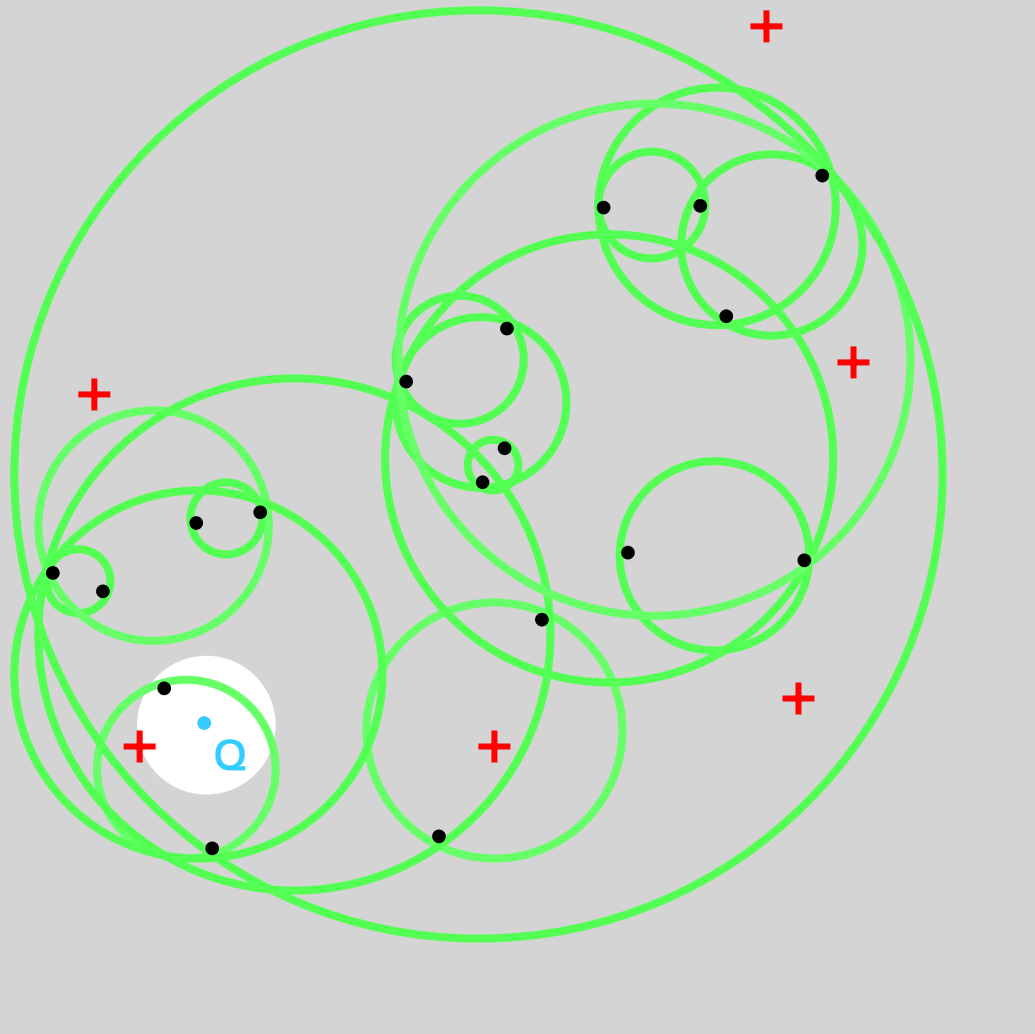




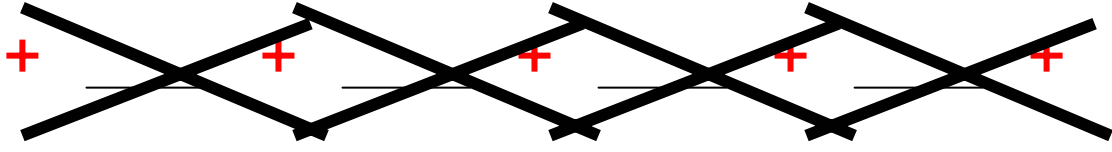
I can't possibly have any interesting points



We're done



\dot{Q} 1



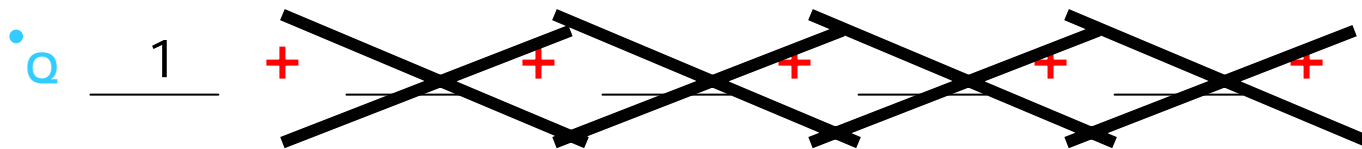
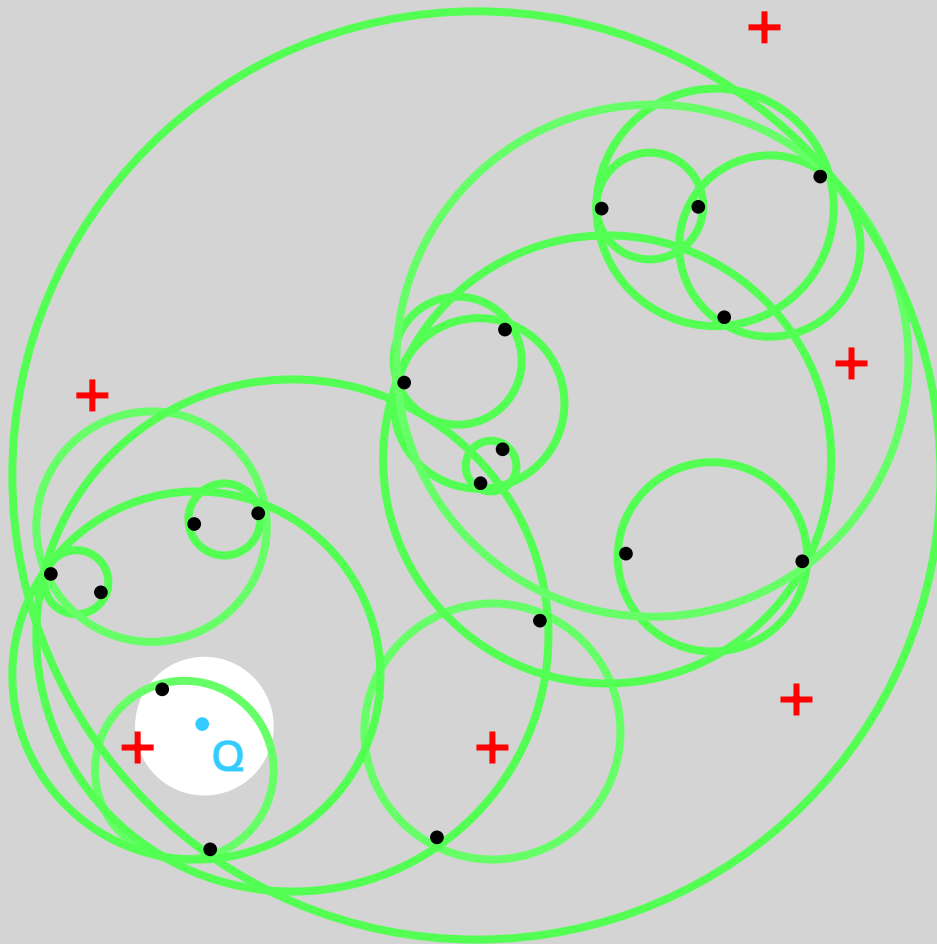
We're done

There's one -ve point closer than the closest +ve point.

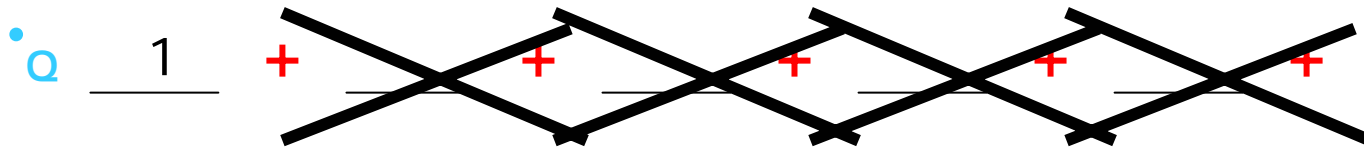
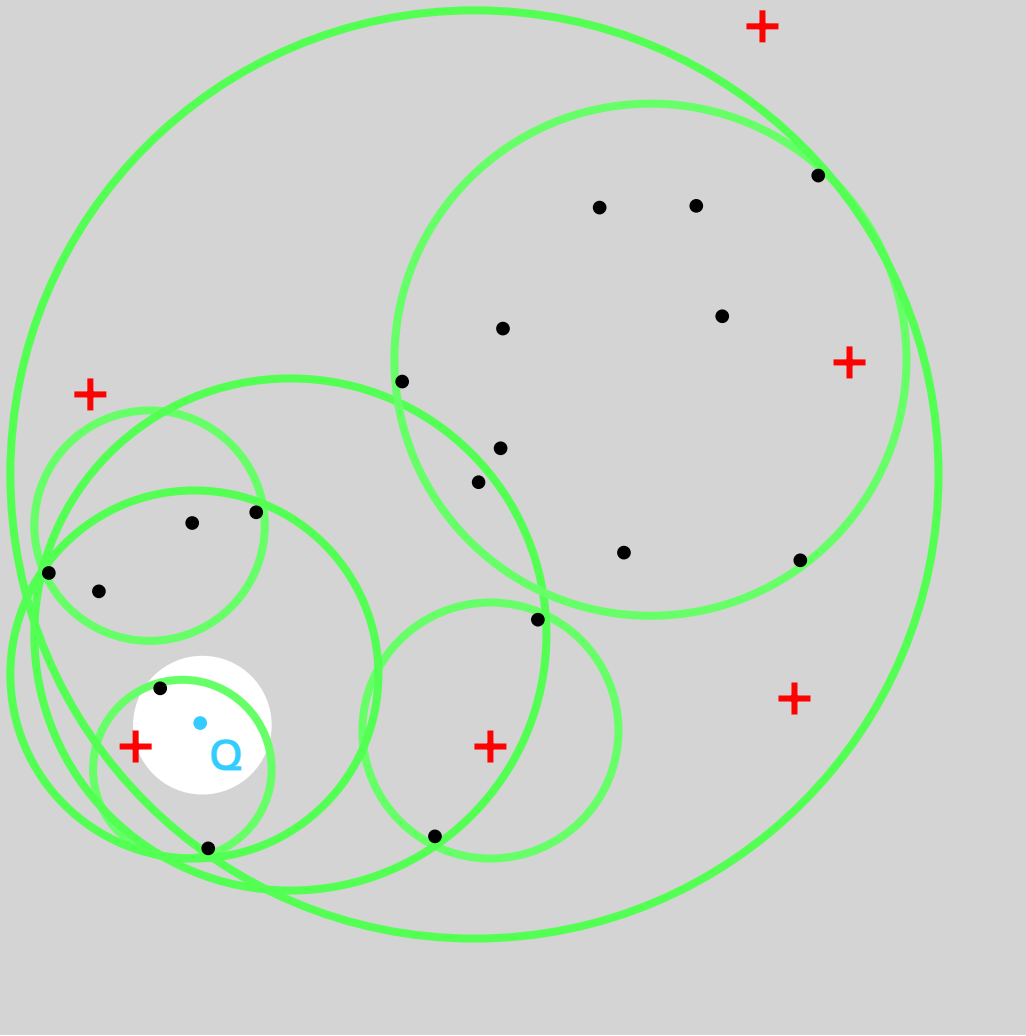
There are more than 3 -ve points closer than the 2nd closest +ve point.

=> Exactly 1 of the 5 nearest neighbors is

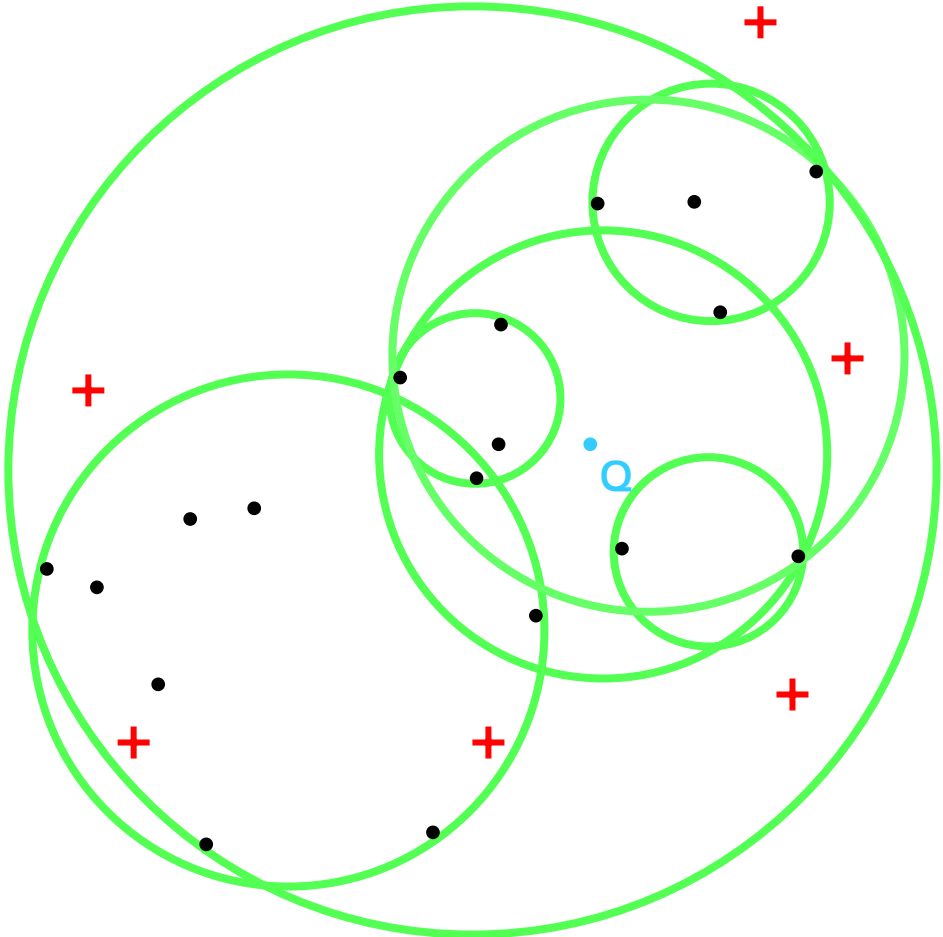
+ve



Balls visited



Another example



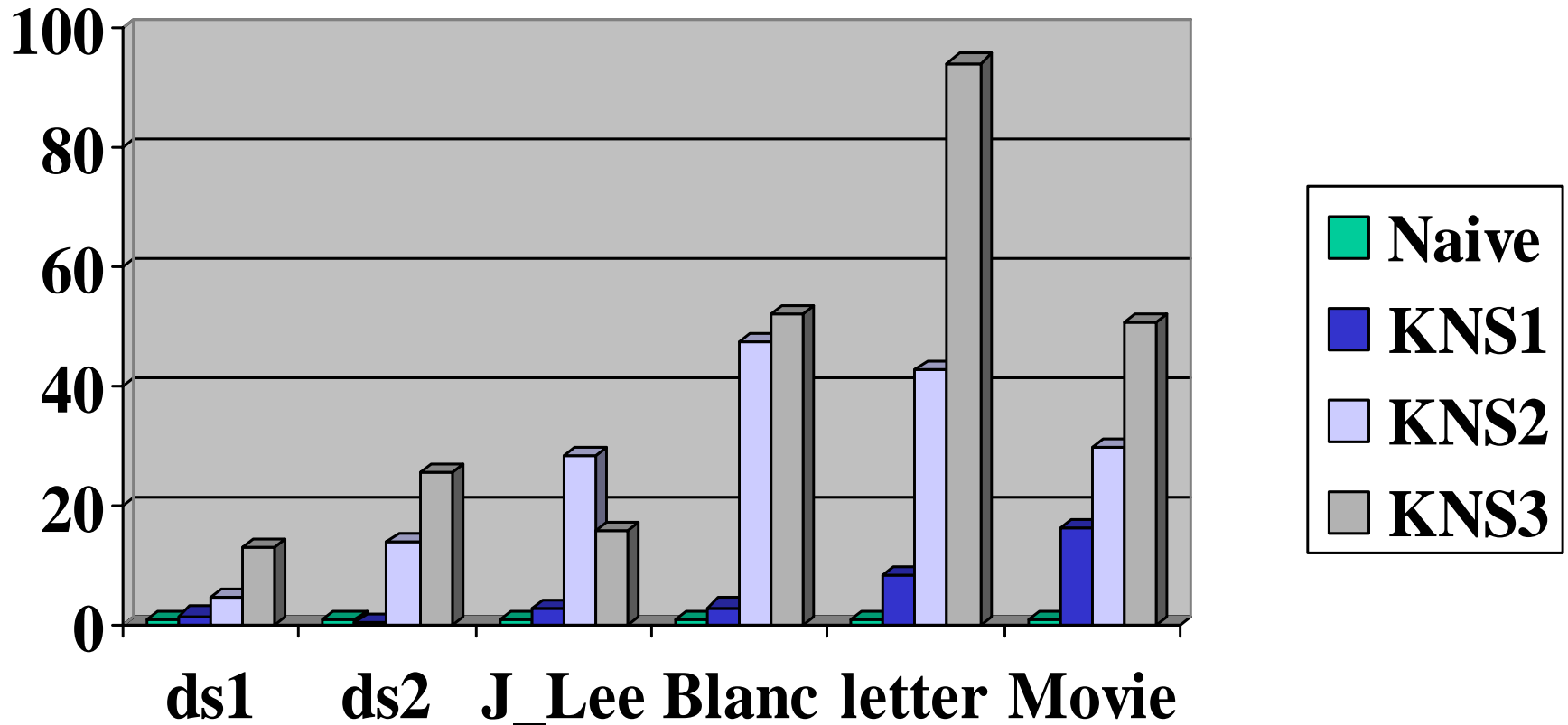
Experimental results

Dataset	Num. of records	Num. of Dimensions	Num.of positive	Num.pos/Num.neg
ds1	26733	6348	804	0.03
ds1.10pca	26733	10	804	0.03
ds1.100pca	26733	100	804	0.03
ds2	88358	1.1×10^6	211	0.002
ds2.100anchor	88358	100	211	0.002
J_Lee.100pca	181395	100	299	0.0017
Blanc_Mel	186414	10	824	0.004

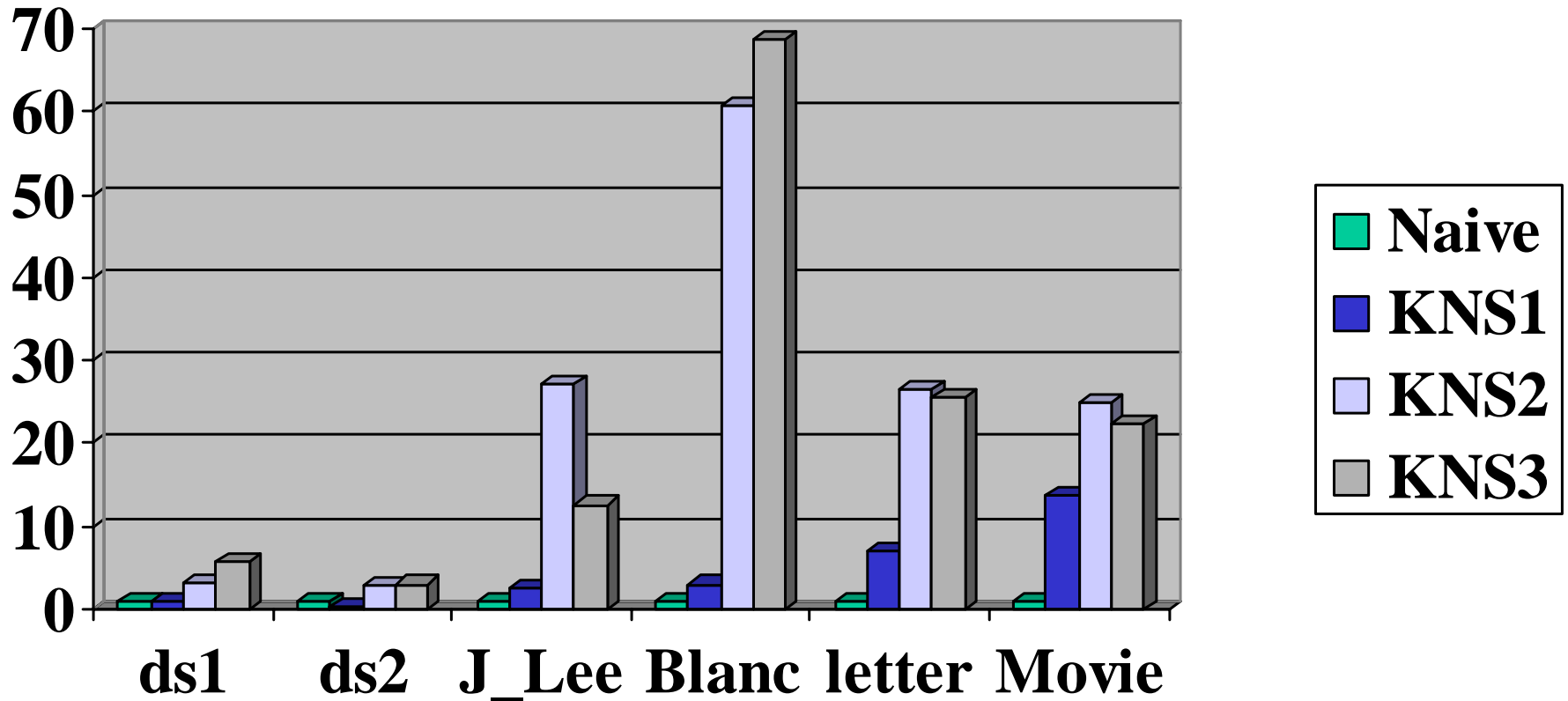
Dataset	Num. records	Num. of Dimensions	Num.of positive	Num.pos/Num.neg
Letter	20000	16	790	0.04
Ipums	70187	60	119	0.0017
Movie	38943	62	7620	0.24
Kdd99(10%)	494021	176	97278	0.24

Num of Distance computations

Speedup for K-NN



Wall-clock-time speedup for k-NN



		NAIVE		KNS1		KNS2		KNS3	
		dists	time (secs)	dists speedup	time speedup	dists speedup	time speedup	dists speedup	time speedup
ideal	k=9	9.0×10^7	30	96.7	56.5	112.9	78.5	4500	486
	k=101			23.0	10.2	24.7	14.7	4500	432
Diag2d(10%)	k=9	9.0×10^7	30	91	51.1	88.2	52.4	282	27.1
	k=101			22.3	8.7	21.3	9.3	167.9	15.9
Diag2d	k=9	9.0×10^9	3440	738	366	664	372	2593	287
	k=101			202.9	104	191	107.5	2062	287
Diag3d	k=9	9.0×10^9	4060	361	184.5	296	184.5	1049	176.5
	k=101			111	56.4	95.6	48.9	585	78.1
Diag10d	k=9	9.0×10^9	6080	7.1	5.3	7.3	5.2	12.7	2.2
	k=101			3.3	2.5	3.1	1.9	6.1	0.7
Noise2d	k=9	9.0×10^7	40	91.8	20.1	79.6	30.1	142	42.7
	k=101			22.3	4	16.7	4.5	94.7	43.5
ds1	k=9	6.4×10^8	4830	1.6	1.0	4.7	3.1	12.8	5.8
	k=101			1.0	0.7	1.6	1.1	10	4.2
ds1.10pca	k=9	6.4×10^8	420	11.8	11.0	33.6	21.4	71	20
	k=101			4.6	3.4	6.5	4.0	40	6.1
ds1.100pca	k=9	6.4×10^8	2190	1.7	1.8	7.6	7.4	23.7	29.6
	k=101			0.97	1.0	1.6	1.6	16.4	6.8
ds2	k=9	8.5×10^9	105500	0.64	0.24	14.0	2.8	25.6	3.0
	k=101			0.61	0.24	2.4	0.83	28.7	3.3
ds2.100-	k=9	7.0×10^9	24210	15.8	14.3	185.3	144	580	311
	k=101			10.9	14.3	23.0	19.4	612	248
J.Lee.100-	k=9	3.6×10^{10}	142000	2.6	2.4	28.4	27.2	15.6	12.6
	k=101			2.2	1.9	12.6	11.6	37.4	27.2
Blanc_Mel	k=9	3.8×10^{10}	44300	3.0	3.0	47.5	60.8	51.9	60.7
	k=101			2.9	3.1	7.1	33	203	134.0
Letter	k=9	3.6×10^8	290	8.5	7.1	42.9	26.4	94.2	25.5
	k=101			3.5	2.6	9.0	5.7	45.9	9.4
Ipums	k=9	4.4×10^9	9520	195	136	665	501	1003	515
	k=101			69.1	50.4	144.6	121	5264	544
Movie	k=9	1.4×10^9	3100	16.1	13.8	29.8	24.8	50.5	22.4
	k=101			9.1	7.7	10.5	8.1	33.3	11.6
Kddcup99 (10%)	k=9	2.7×10^{11}	1670000	4.2	4.2	574	702	4	4.1
	k=101			4.2	4.2	187.7	226.2	3.9	3.9

Outline

Cached Sufficient Statistics

Ball Trees (= Metric Trees)

K-nearest neighbor with ball trees

Very fast non-parametric classification

skewed binary outputs

General binary outputs

multi-classed outputs

▶ Very fast kernel-based statistics

n-point computations

clustering

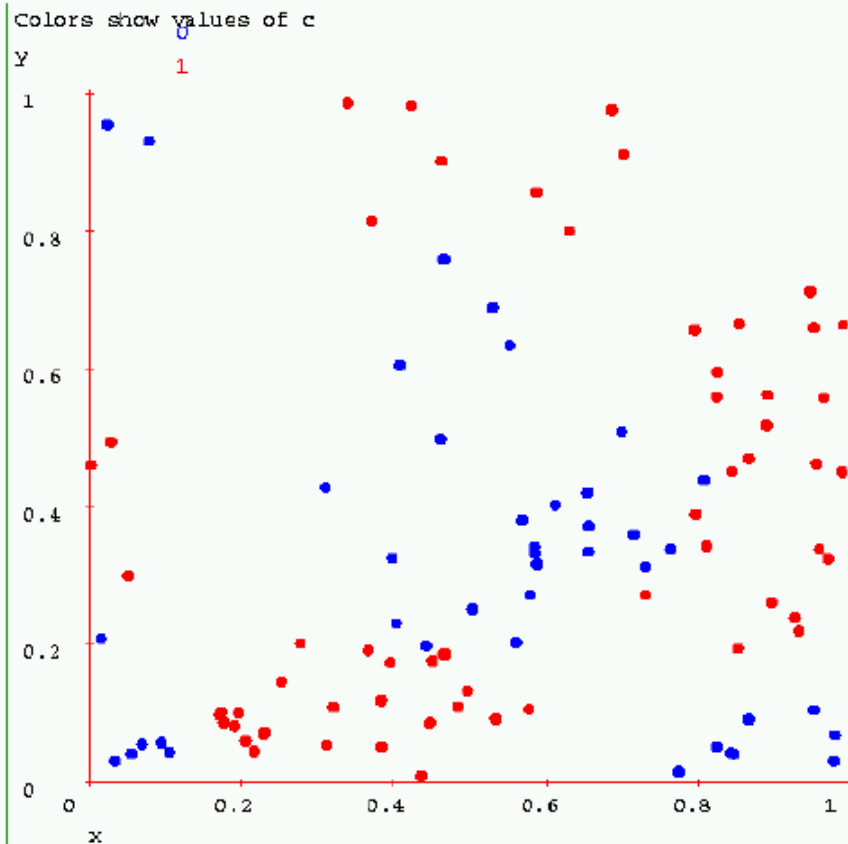
non-parametric clustering (overdensity hunting)

Active learning for anomaly hunting

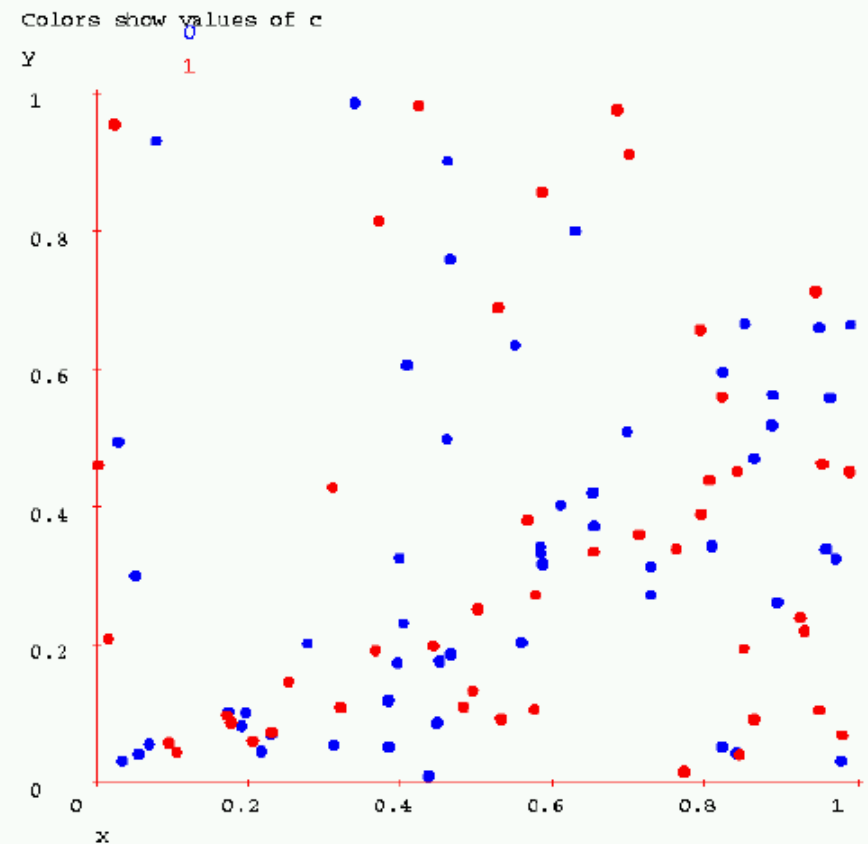
GMorph: Efficient Galaxy morphology fitting

Other Auton topics

All-pairs-of-points problems in statistics



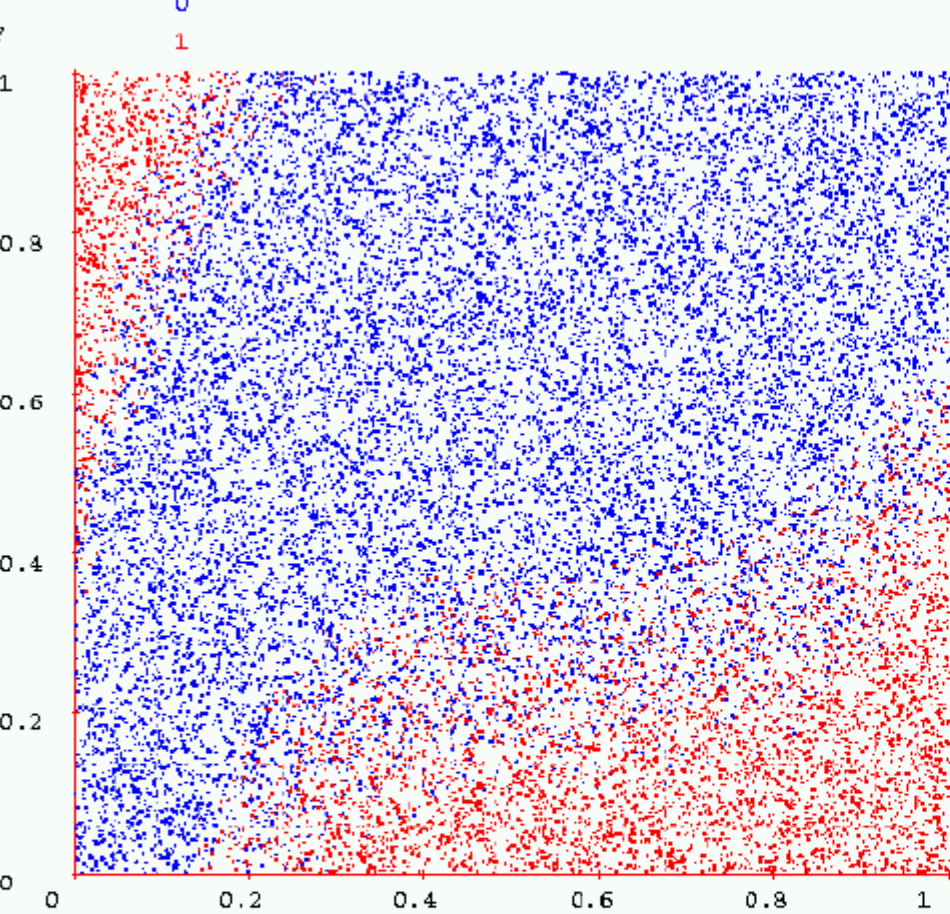
Leave-one-out testing with one nearest neighbor says
“I can see something interesting”



Leave-one-out testing with one nearest neighbor says
“I can see nothing interesting”

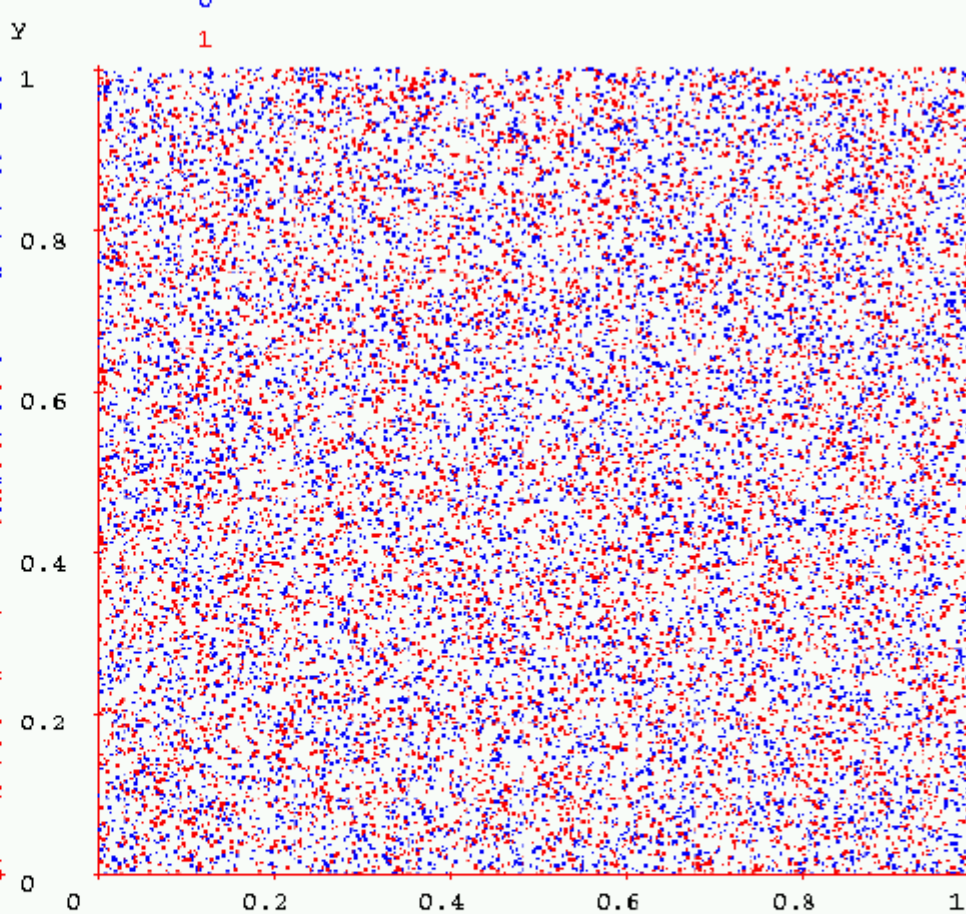
Kernel Regression

Colors show values of c



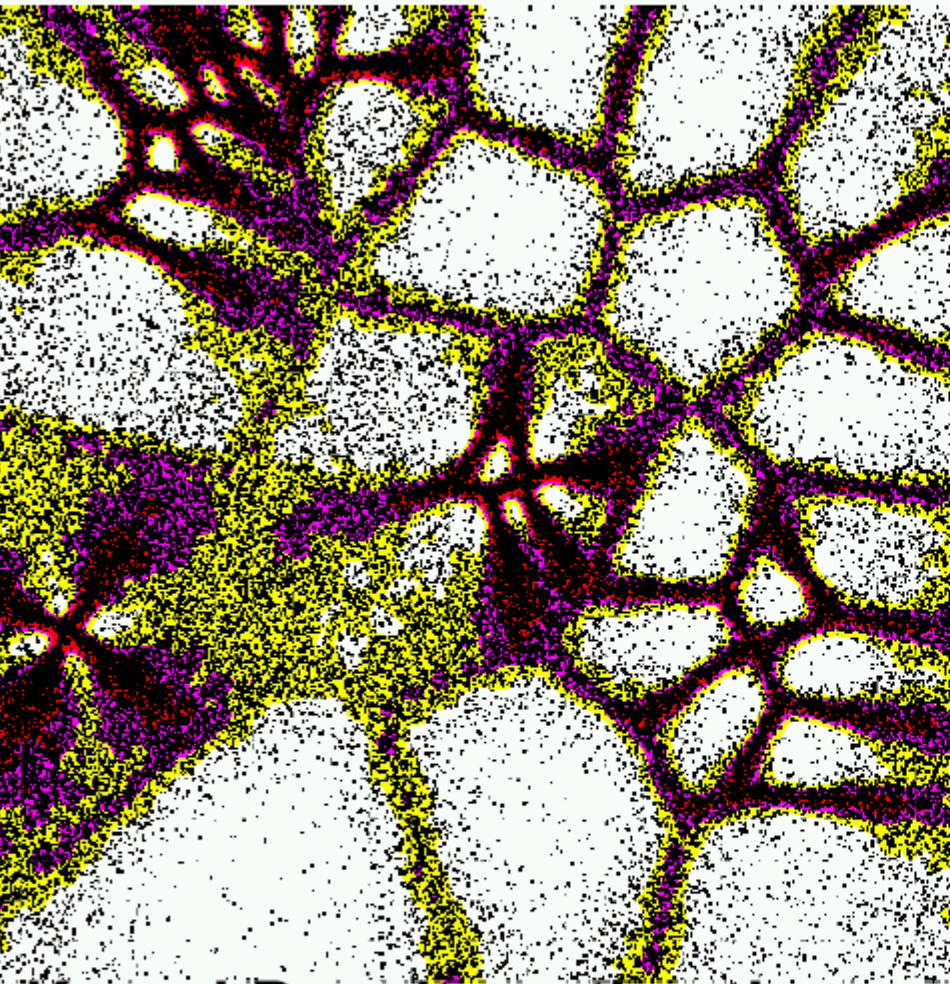
Kernel regression says “I can see something interesting”

Colors show values of k

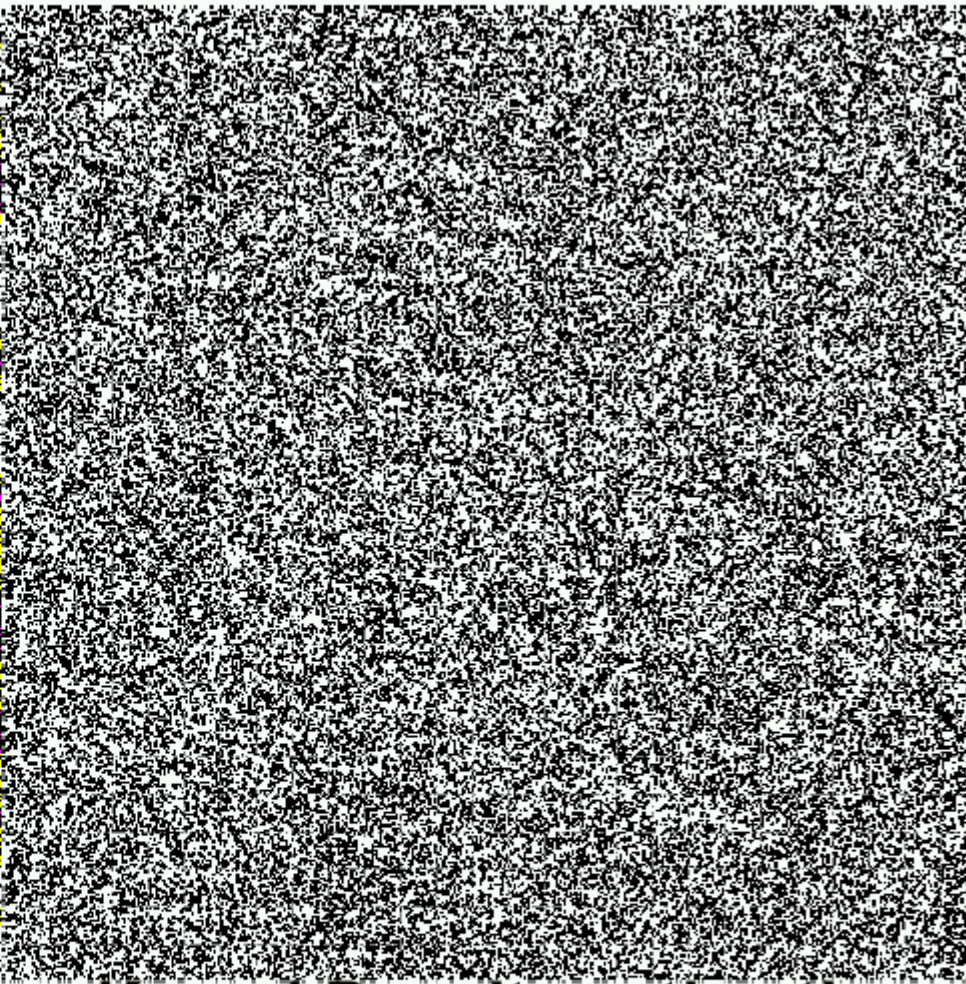


Kernel Regression says “I can see nothing interesting”

Kernel Density Estimation

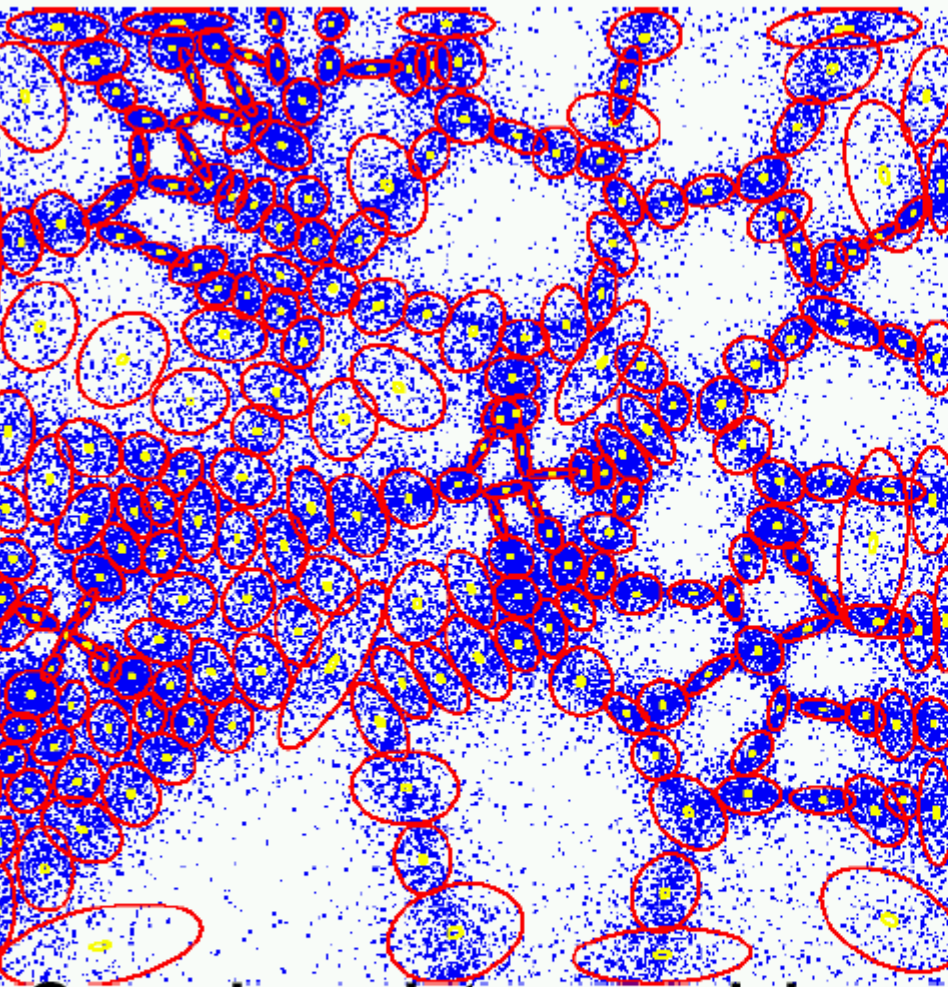


Kernel Density estimation
say “I can see something
interesting”

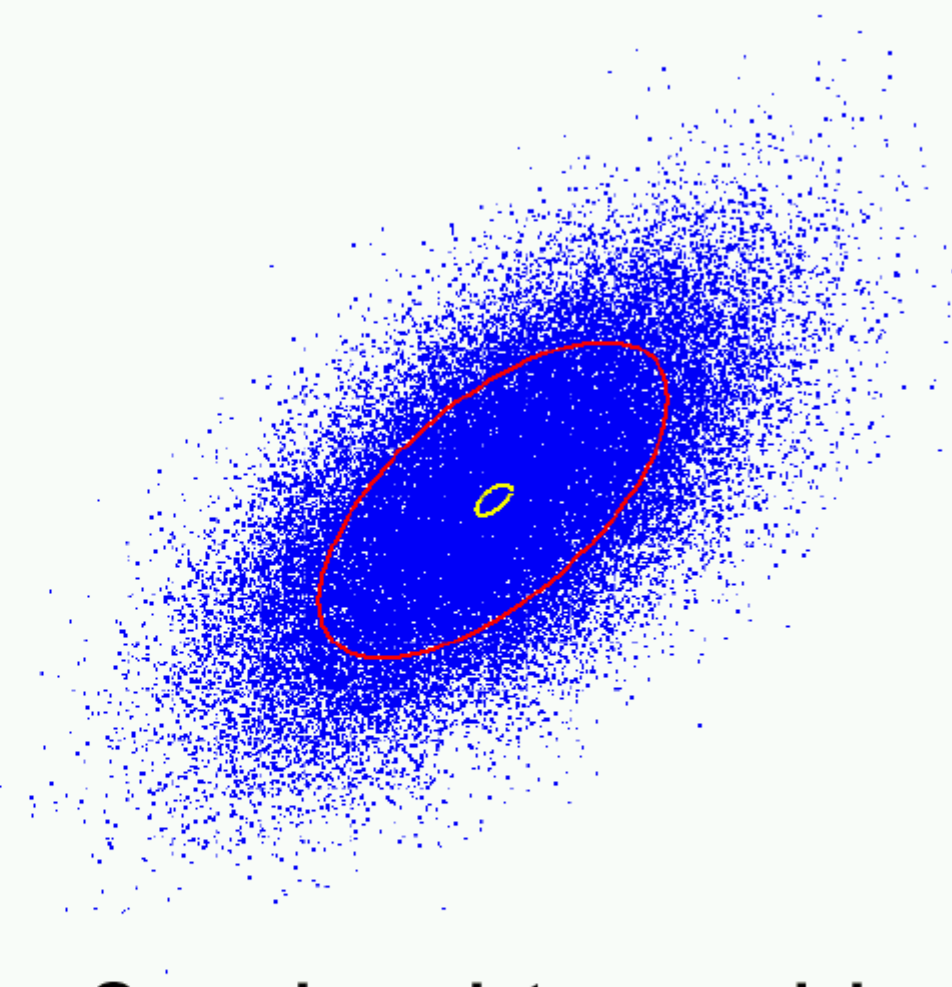


Kernel Density estimation
says “I can see nothing
interesting”

Many-component Mixture Models



Gaussian mixture model
says “I can see and
estimate a great deal of
structure”

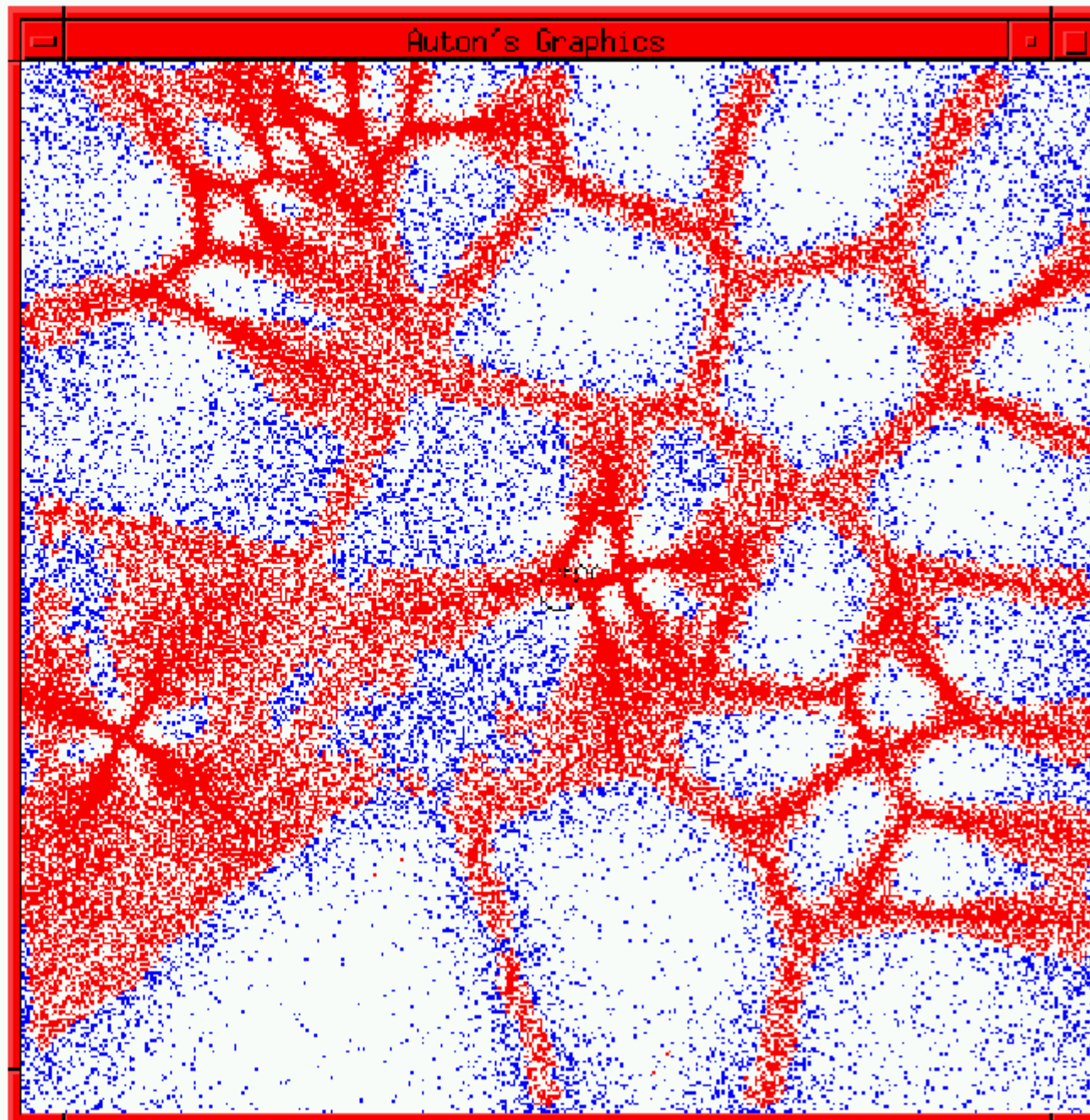


Gaussian mixture model
says “It’s a big old lump”

Spatial Anomaly Detection

Red dots are in a crowded neighborhood.

Blue dots are lonely.



Many other “All-pairs” problems

- Locally weighted polynomial regression
- Gaussian processes
- Point processes
- Bottom-up clustering

“All-pairs of attributes” important too...

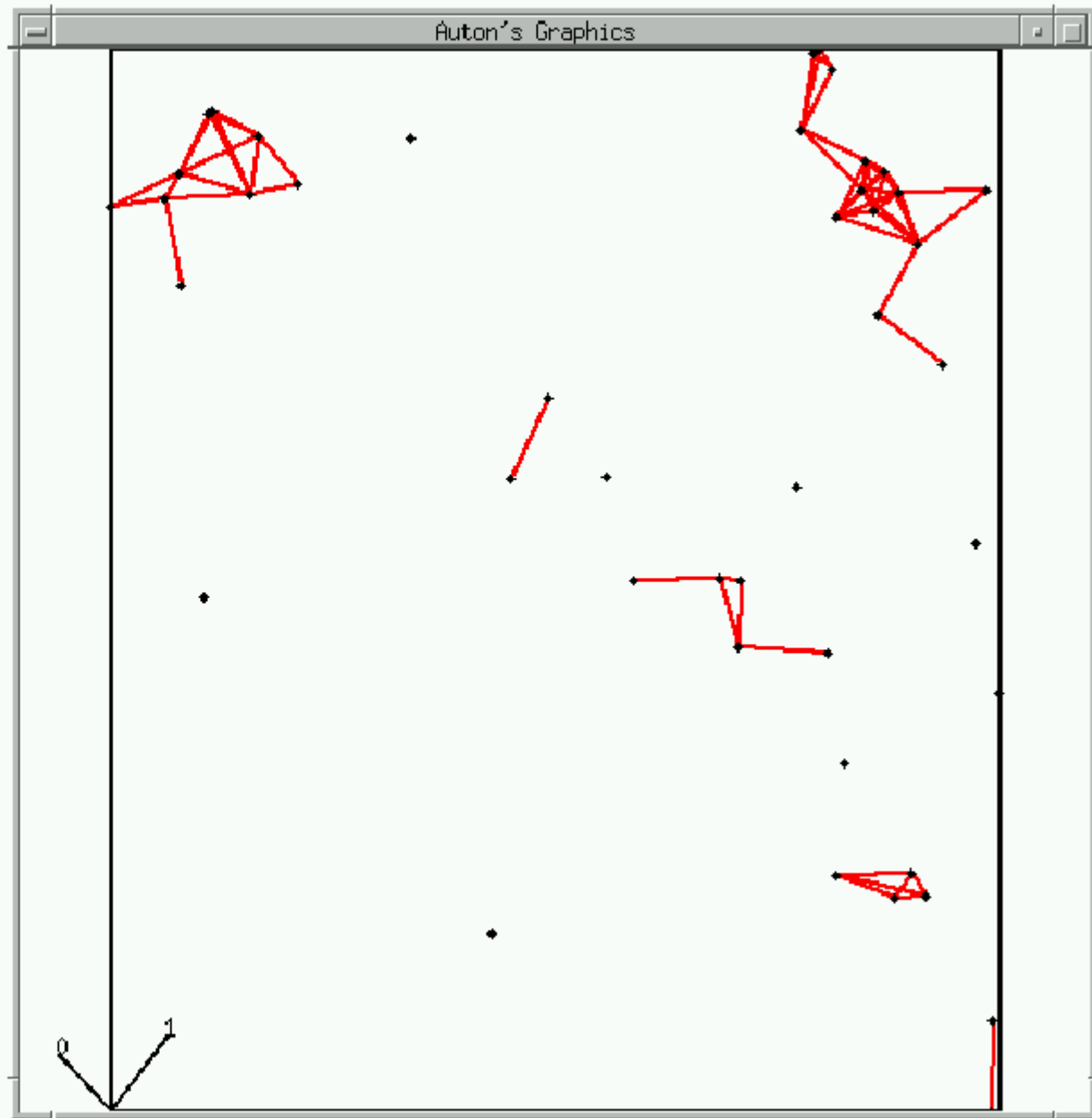
- Find me the most highly correlated pair of attributes.
- The most similar color-bands, image filters...

2-point correlation

...the purest form of an “all-pairs” problem.

There are 62 pairs of points that lie within 0.1 units of each other:

..important in astrophysics for characterizing matter distribution.



Fast all-point-pairs: Idea One

Use an $O(n^2)$ algorithm and buy a fast computer

Problem: $O(n^2)$ is vicious.

Comparative Results

Non-approximate version

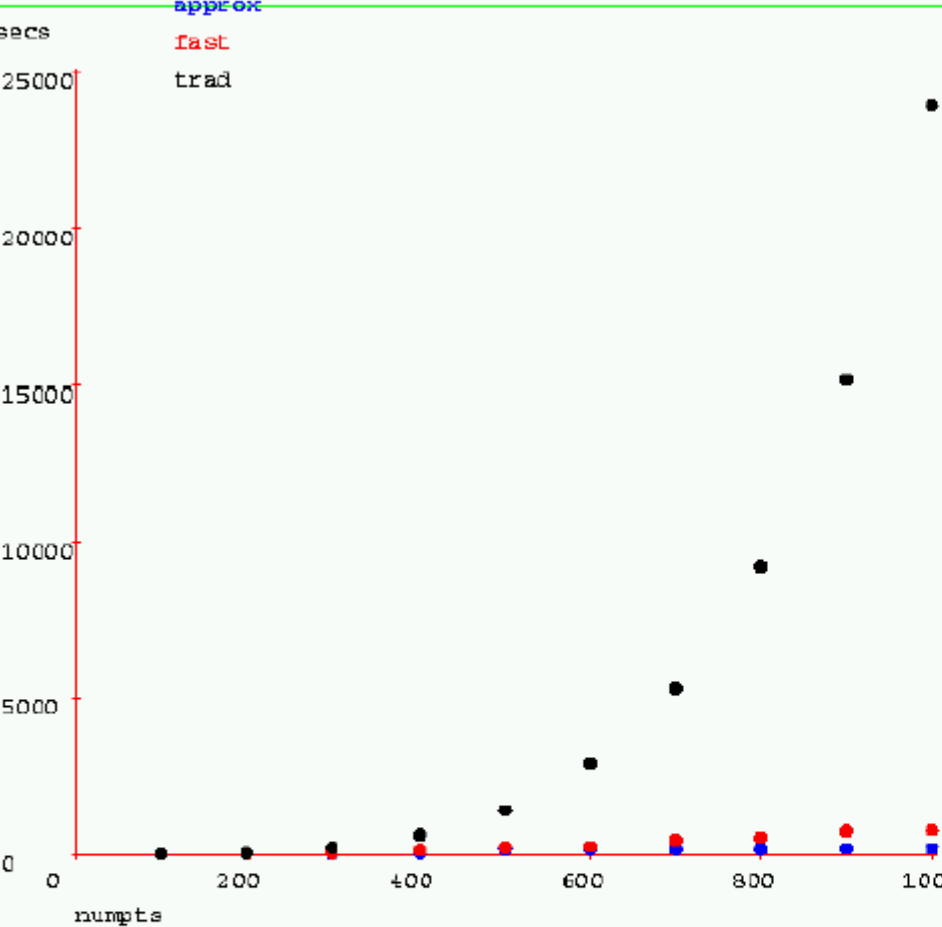
Number of Points	Quadratic time (secs)	Single-tree time (secs)	Dual-tree time (secs)	Single Tree Speedup	Dual Tree Speedup
10000	132	2.2	1.2	60	110
20000	528	4.8	2.8	110	189
50000	3300	11.8	7.0	280	471
150000	30899	37	20	835	1545
300000	123599	76	40	1626	3090

Approximate version (20,000 datapoints on slower machine):

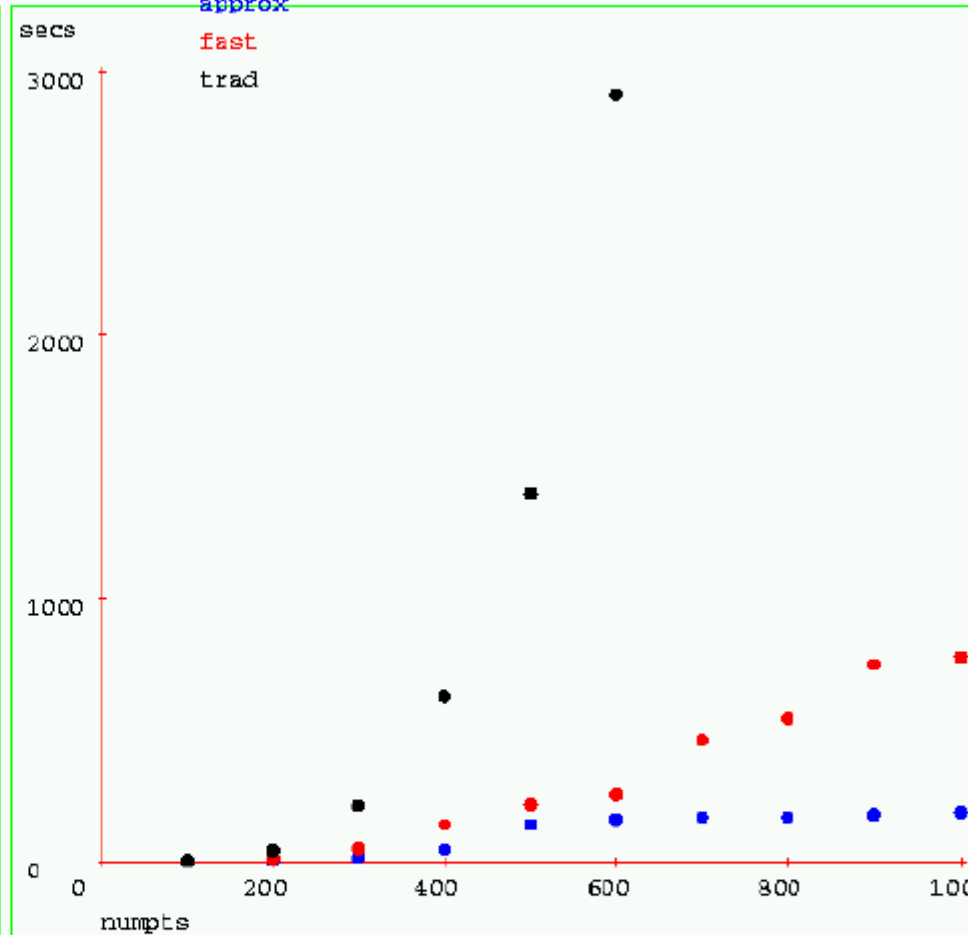
ϵ	0.001	0.01	0.02	0.05	0.1	0.2	0.5
secs	37	30	30	18	10	10	0.3

4-point performance

Colors show values of method



Colors show values of method



Black: Traditional 4-point
Red: Fast Exact 4-point
Blue: Fast Approx 4-pt

Closeup

Outline

Cached Sufficient Statistics

Ball Trees (= Metric Trees)

K-nearest neighbor with ball trees

- Very fast non-parametric classification

- skewed binary outputs

- General binary outputs

- multi-classed outputs

Very fast kernel-based statistics

- n-point computations



- clustering

- non-parametric clustering (overdensity hunting)

Active learning for anomaly hunting

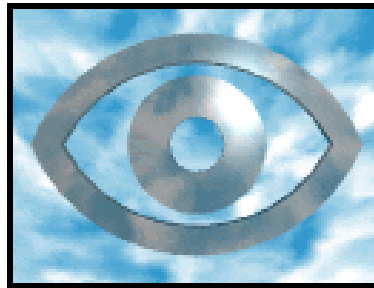
GMorph: Efficient Galaxy morphology fitting

Other Auton topics

Data Structures for Fast K-means

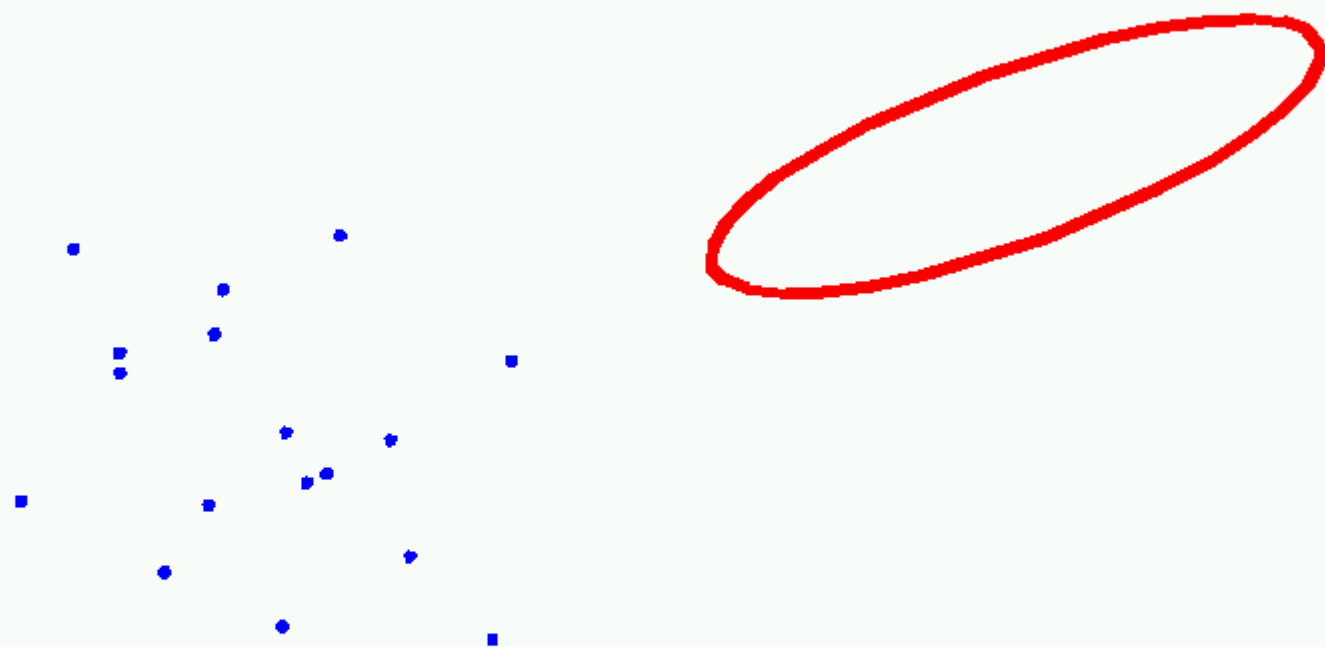
The Auton Lab

Carnegie Mellon University



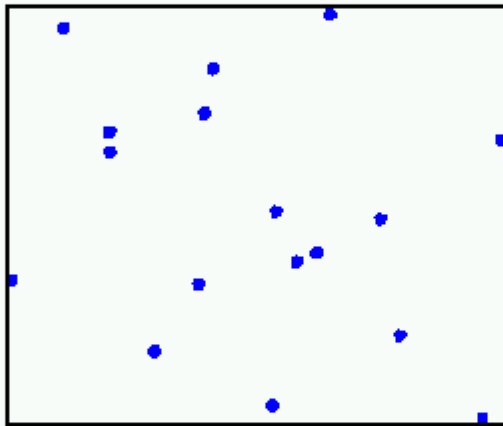
www.autonlab.org

Computing likelihood of datapoints...



Suppose you want to compute the sum of log-likelihoods of all the blue dots given they'd been generated by the big red Gaussian.

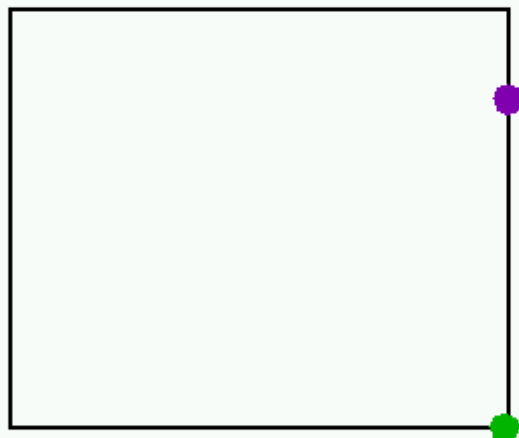
Computing likelihood of datapoints...



Suppose we happen to know their bounding box

Computing likelihood of datapoints...

**Greatest log-likelihood
if they're all here**

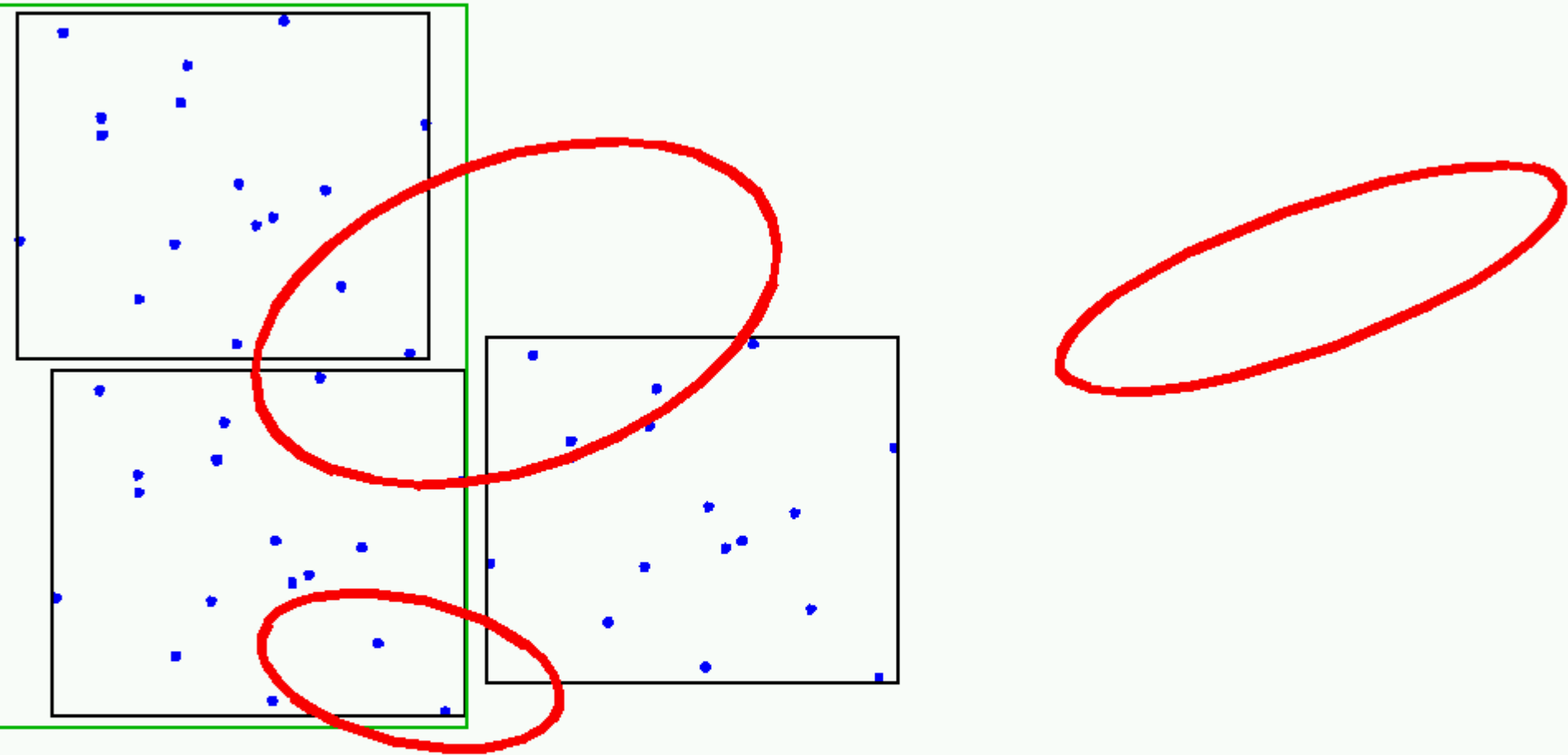


**Least log-likelihood
if they're all here**



**Without visiting the points individually, we can put bounds
on their contributions to the Gaussian. Sometimes those
bounds'll be tight enough...**

Many points, boxes, Gaussians...



If you play this game on a large scale, you find yourself doing
kernel densities, locally weighted regression, locally
weighted PCA, adaptive kernels, k-means clustering, snake-
based filament tracing, hierarchical classification.....
.....very very very fast

Cached sufficient statistics

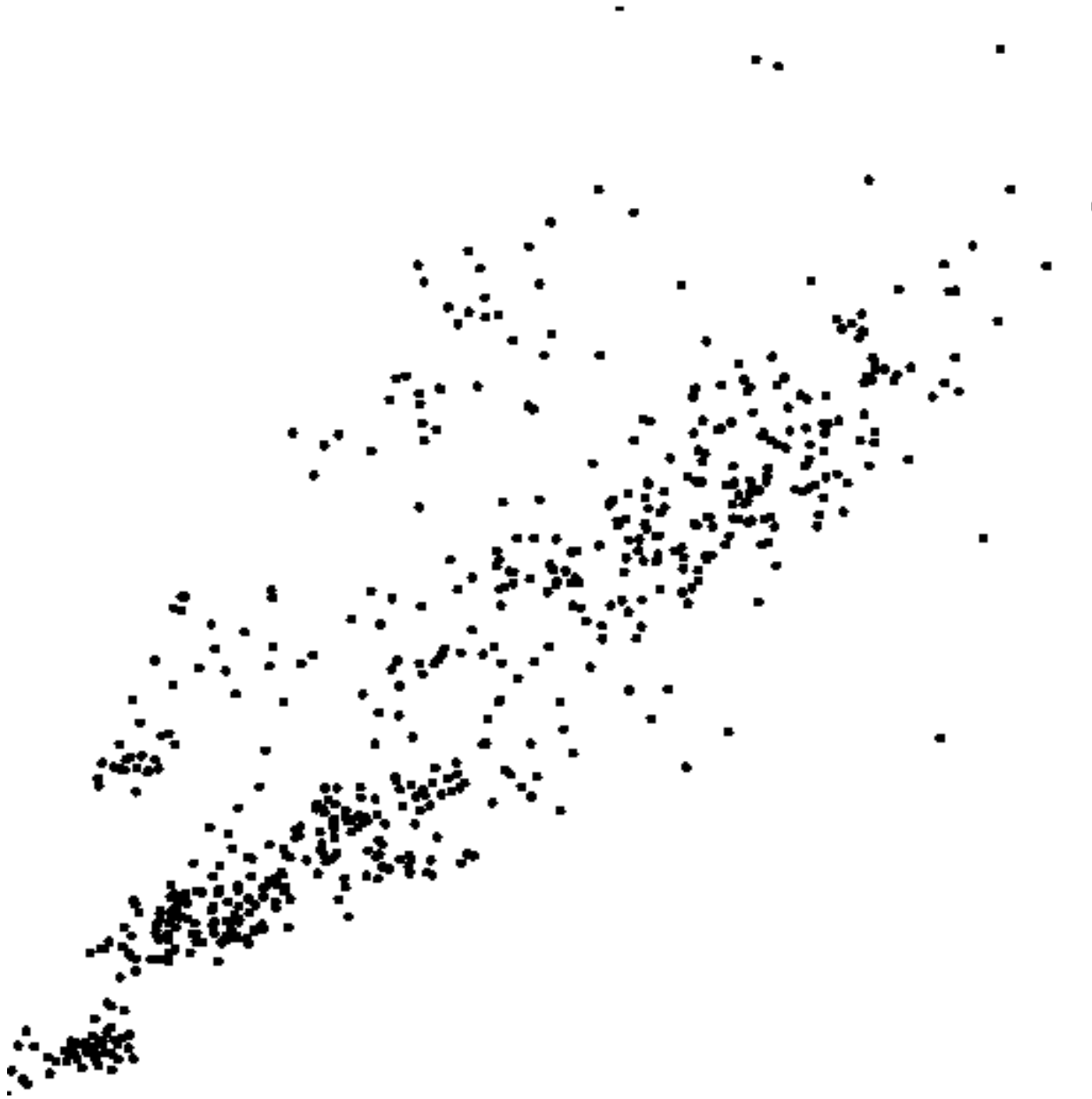
What I've shown you:

- It's intuitively possible to look at a node in a tree and decide whether in order to estimate the data loglikelihood you need to see more detail.

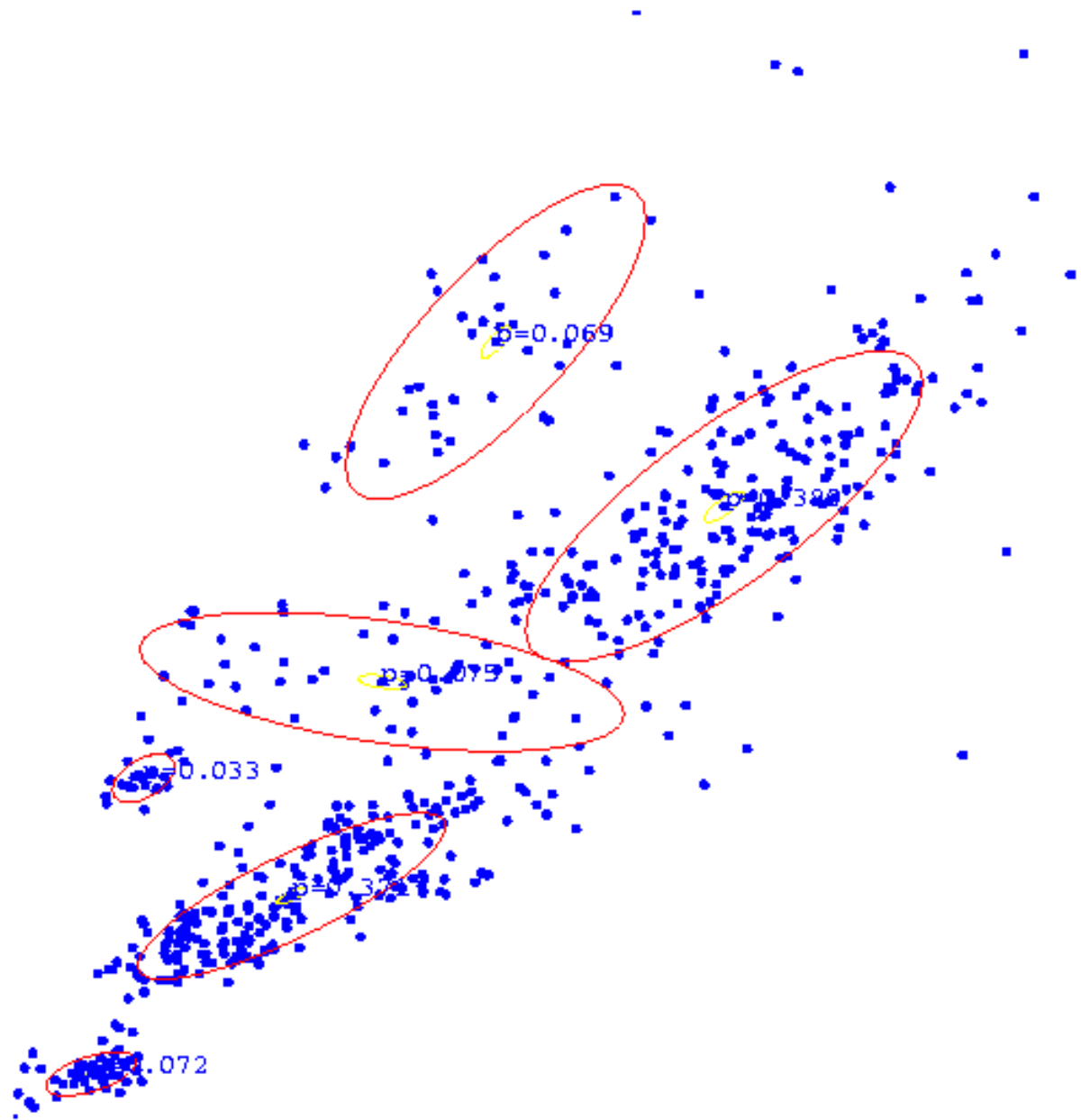
What I don't have time to show you

- Why you must also cache information other than the bounding box in every single kd-tree node:
 - The centroid of all points it owns.
 - The covariance of all points.
- Each algorithm plays different tricks with these kinds of bounds
- Same principal as Barnes Hut and Greengard but sometimes trickier.

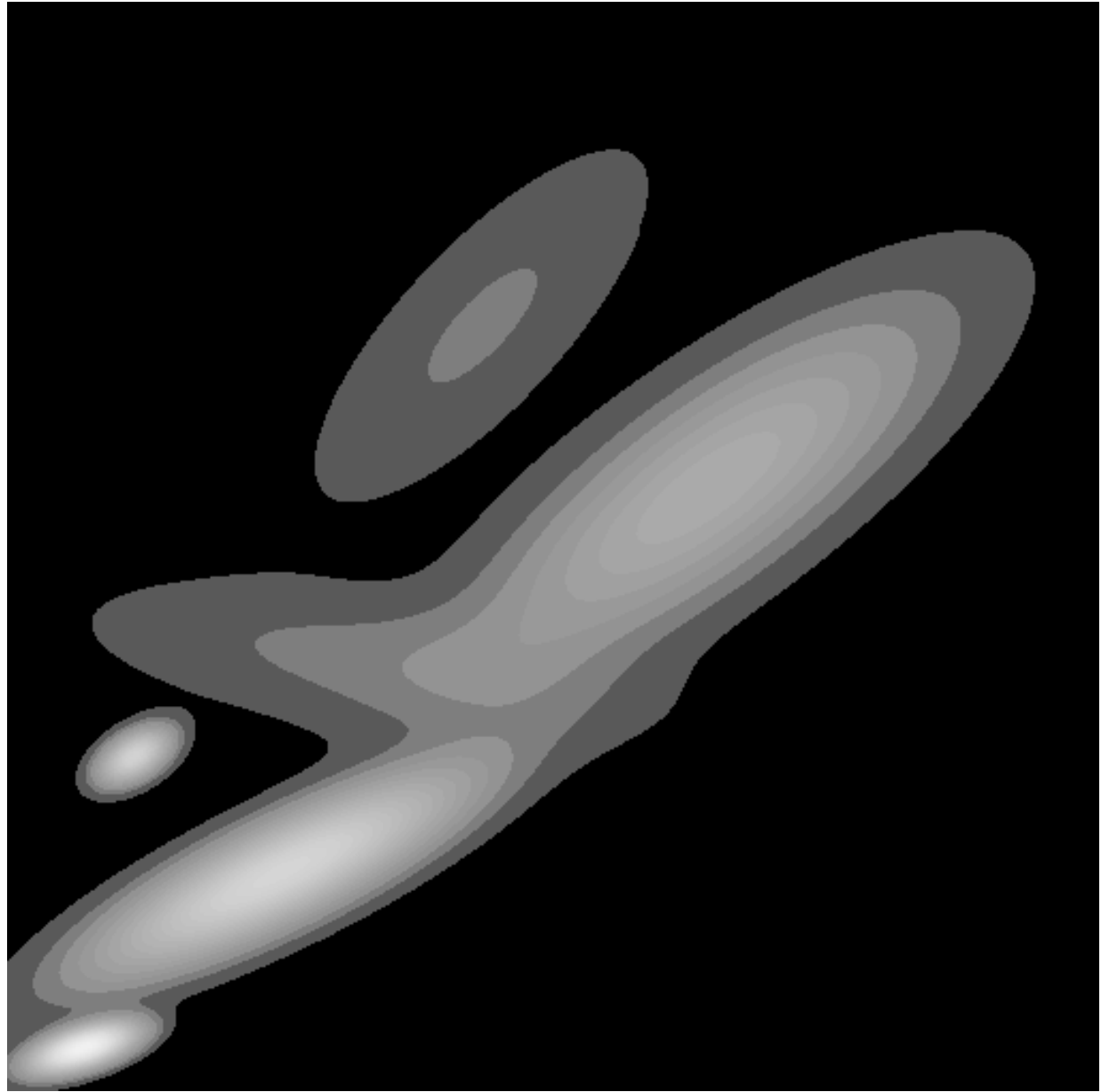
Some Bio Assay data



GMM clustering of the assay data



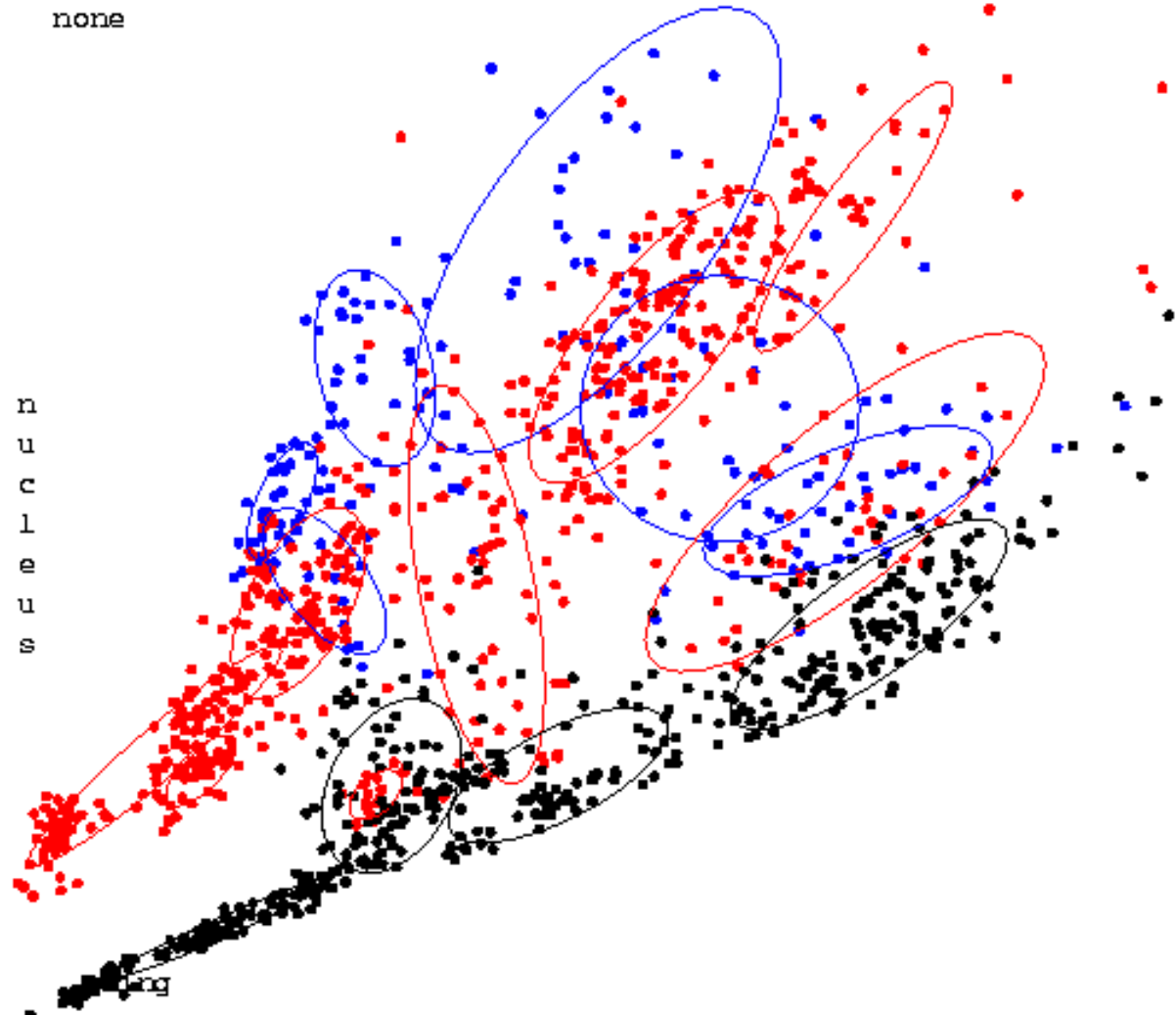
Resulting Density Estimator



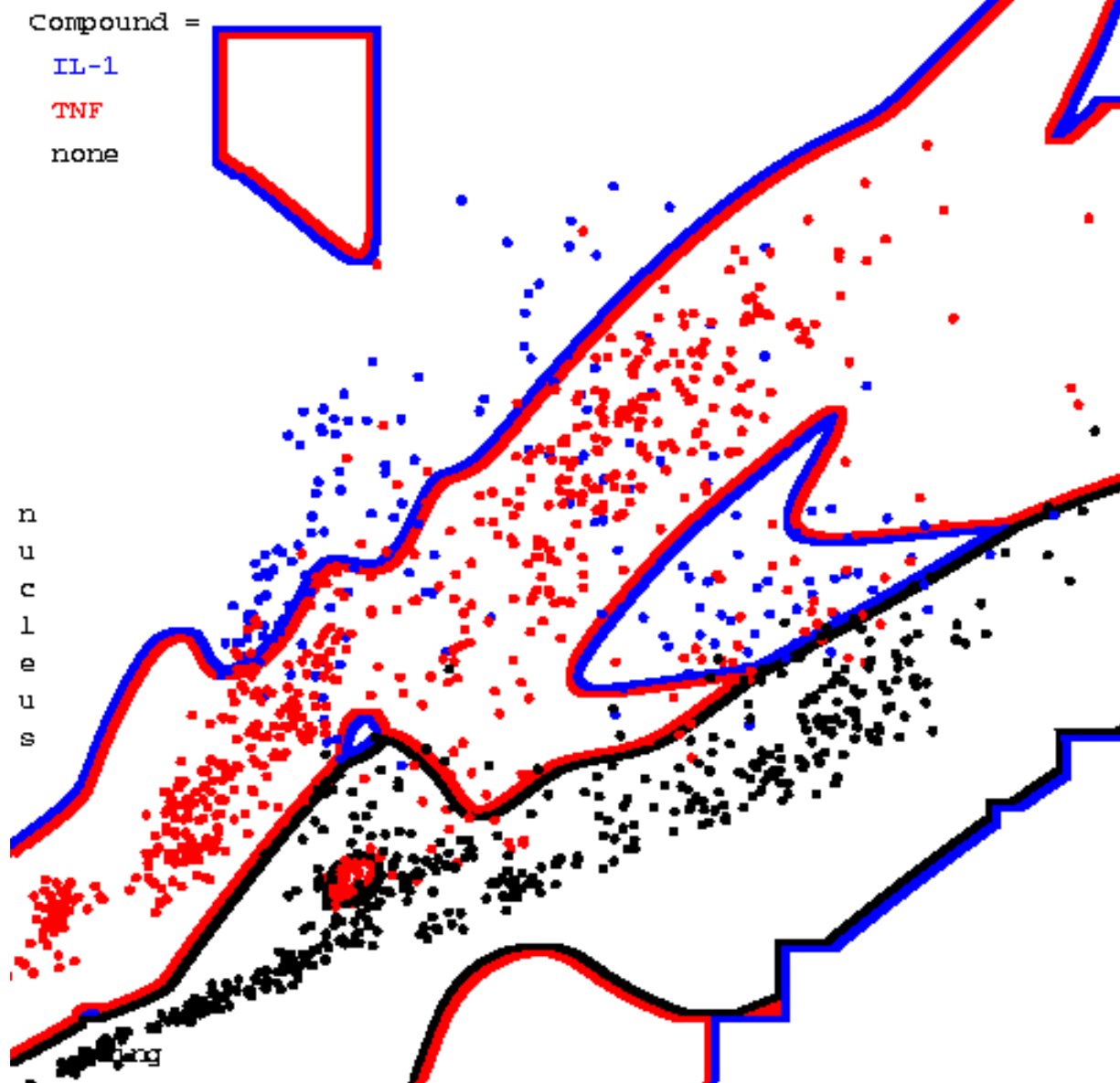
Three classes of assay

(each learned with its own mixture model)
(Sorry, this will again be semi-useless in black and white)

Compound =
IL-1
TNF
none



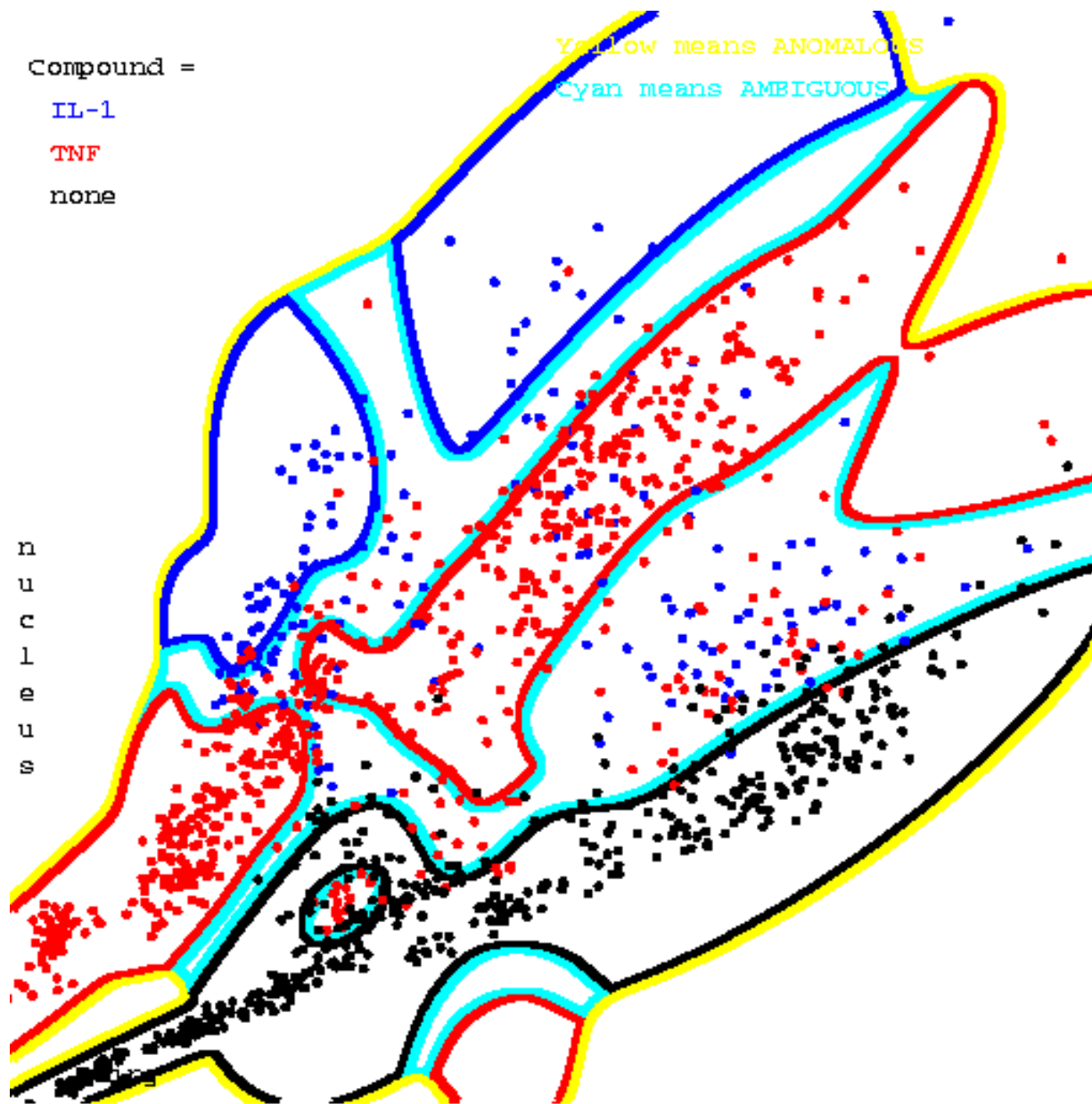
Resulting Bayes Classifier



Resulting Bayes Classifier, using posterior probabilities to alert about ambiguity and anomalousness

Yellow means anomalous

Cyan means ambiguous



Outline

Cached Sufficient Statistics

Ball Trees (= Metric Trees)

K-nearest neighbor with ball trees

Very fast non-parametric classification

skewed binary outputs

General binary outputs

multi-classed outputs

Very fast kernel-based statistics

n-point computations

clustering

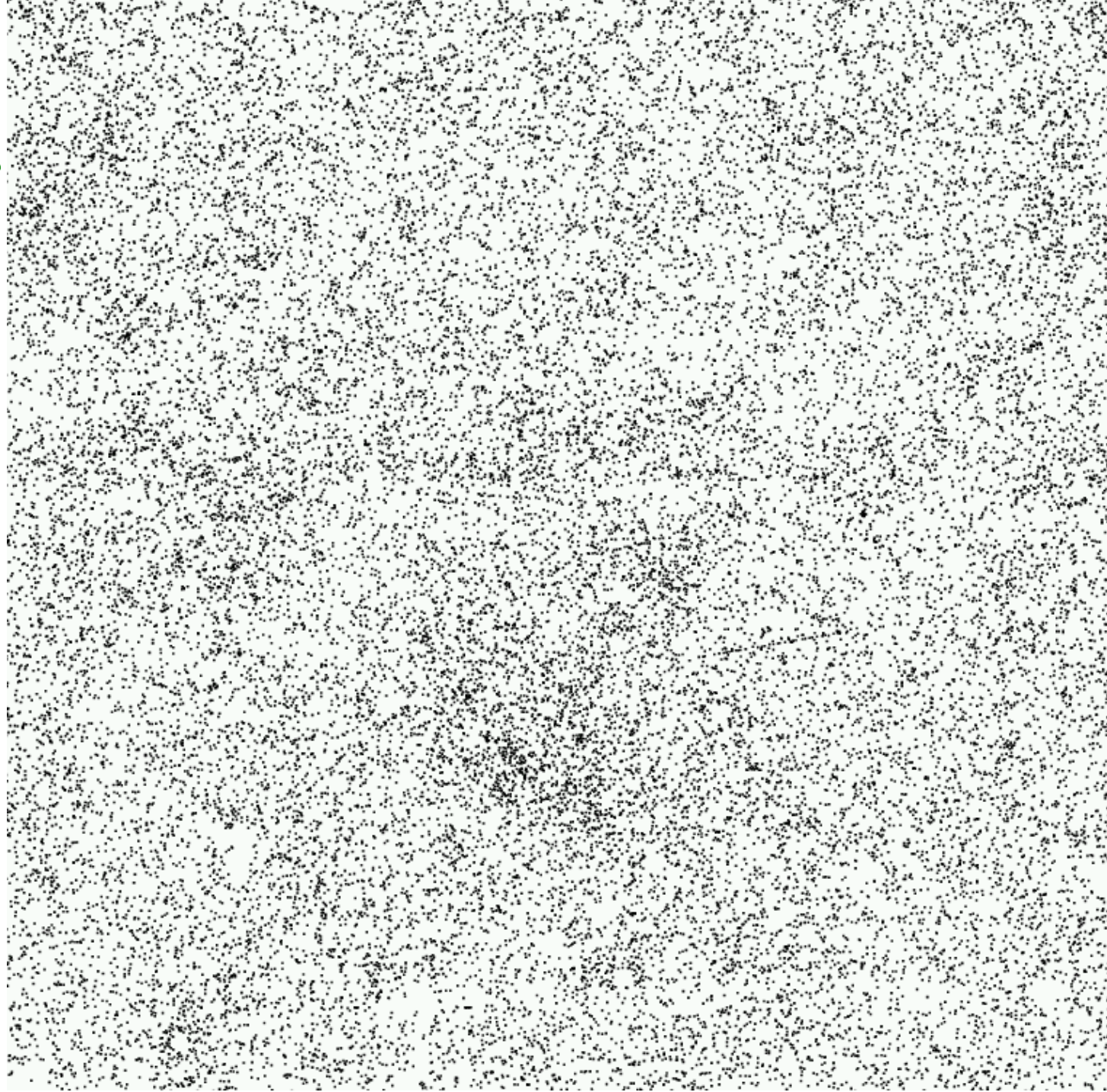
▶ non-parametric clustering (overdensity hunting)

Active learning for anomaly hunting

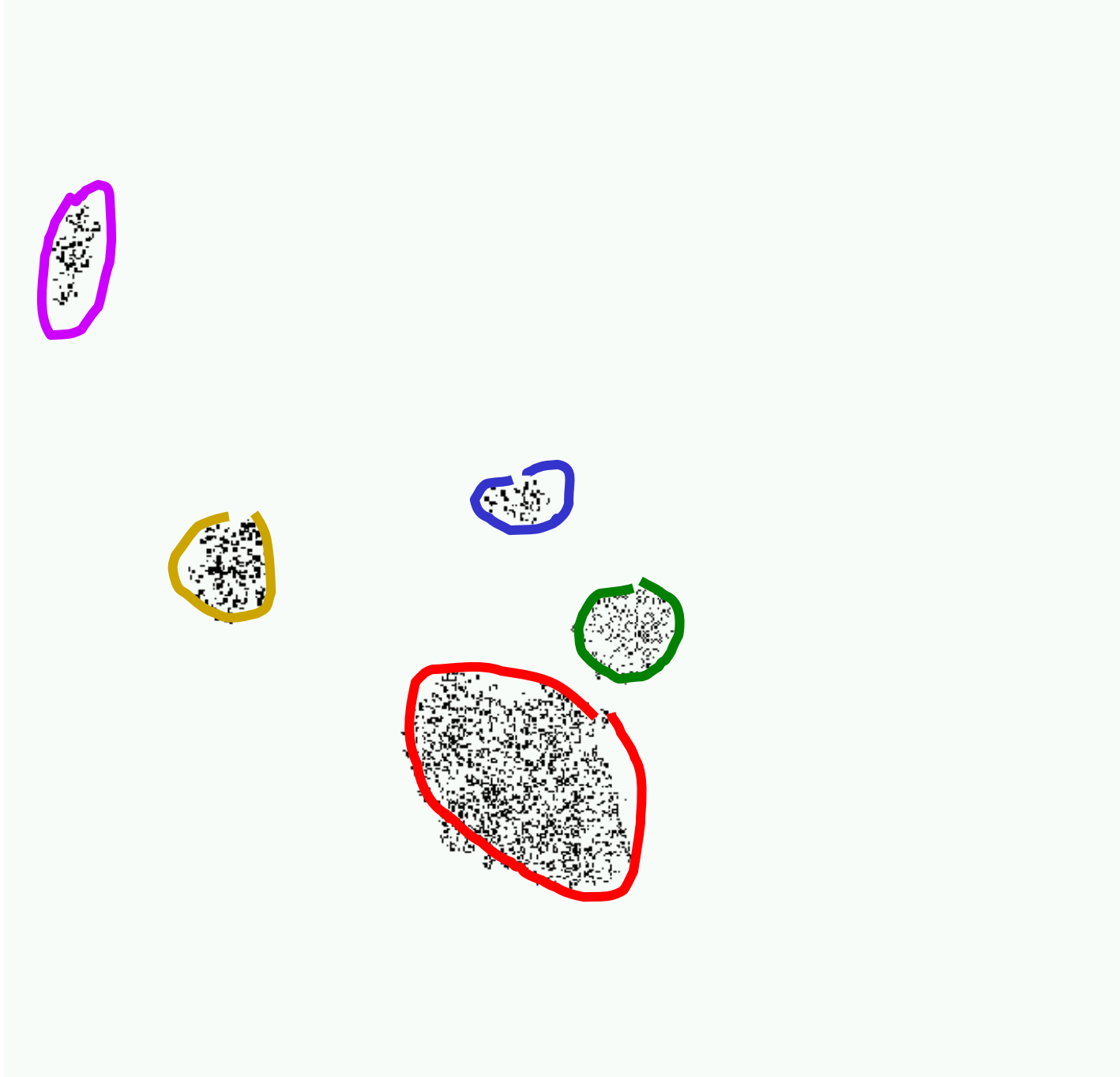
GMorph: Efficient Galaxy morphology fitting

Other Auton topics

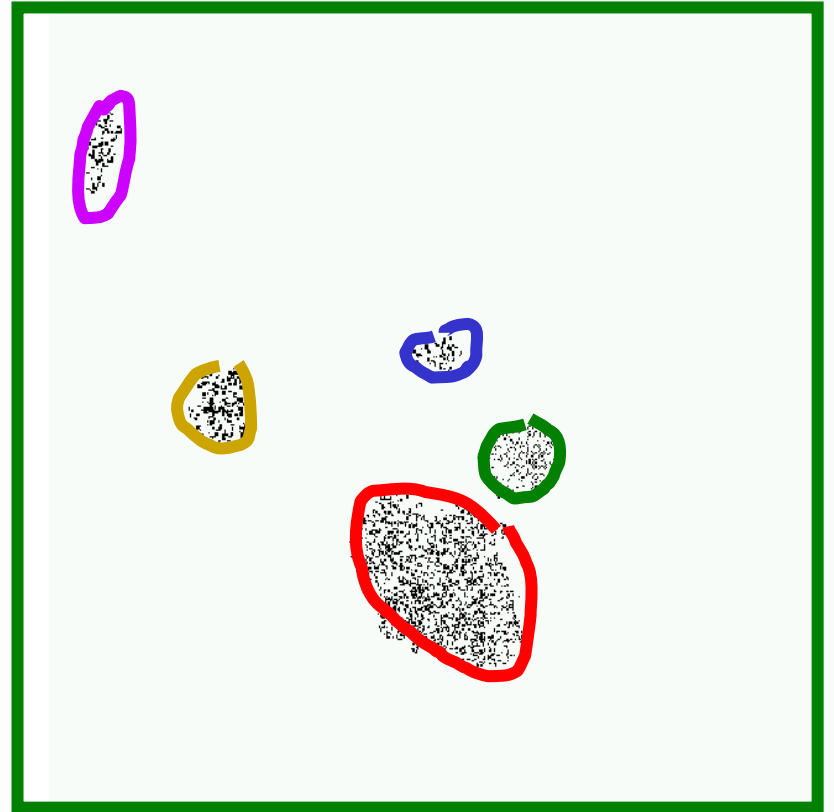
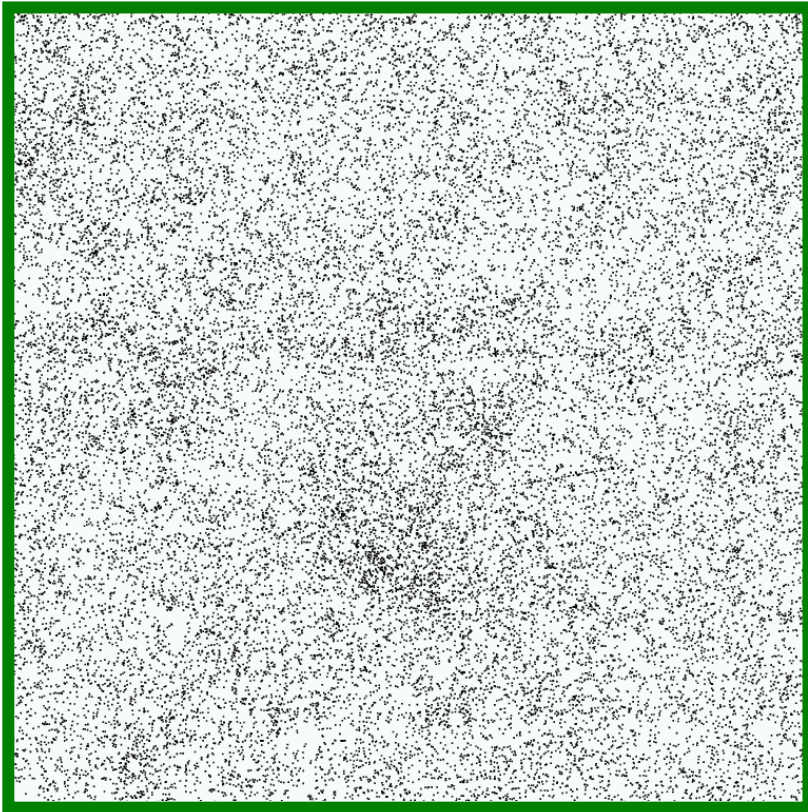
Detecting overdensities



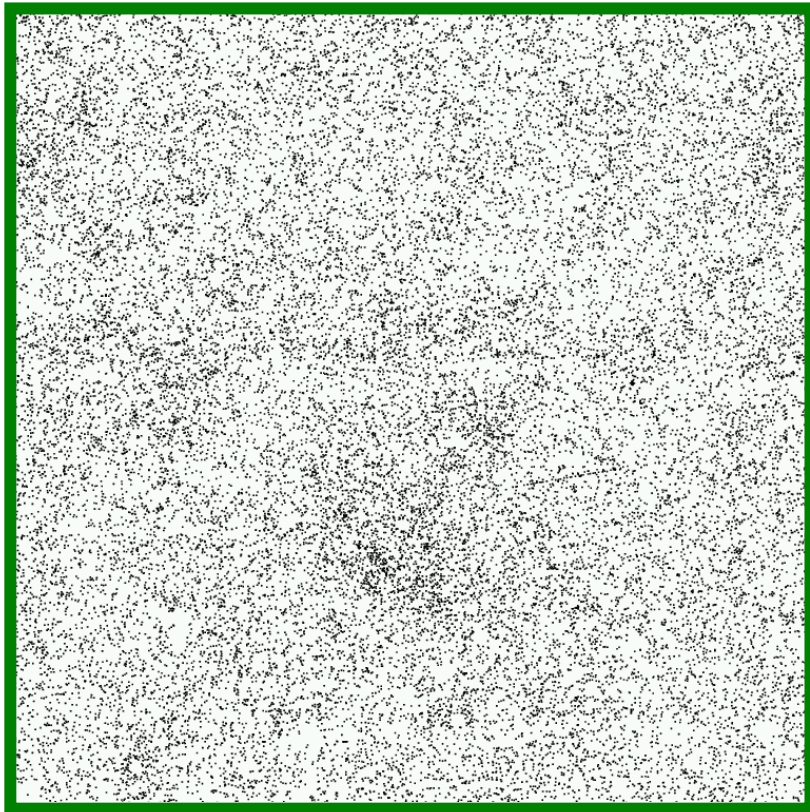
Detecting overdensities



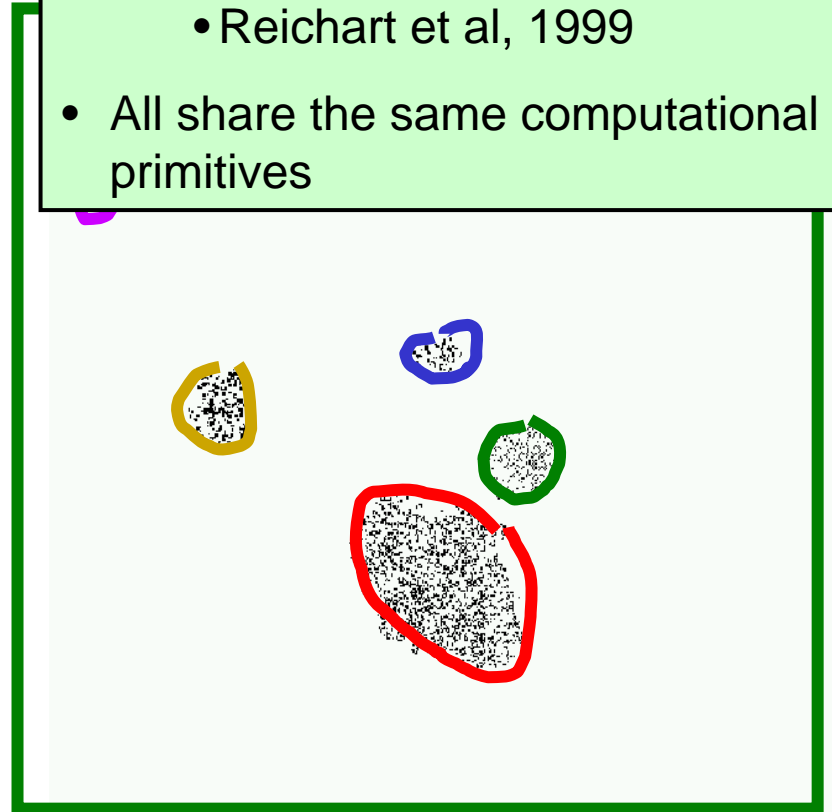
Finding the overdense regions



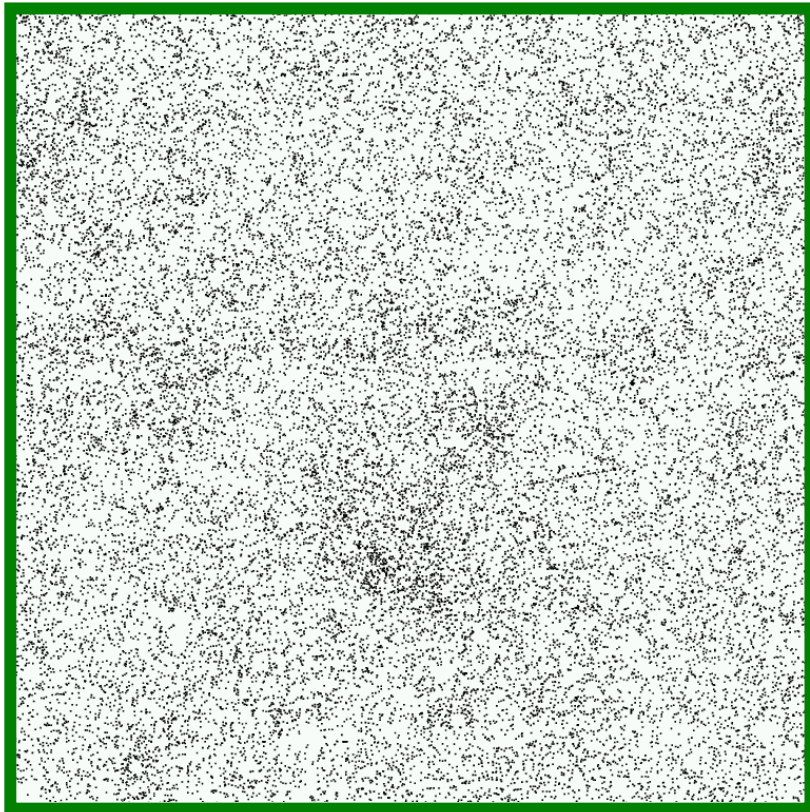
Finding the overdense regions



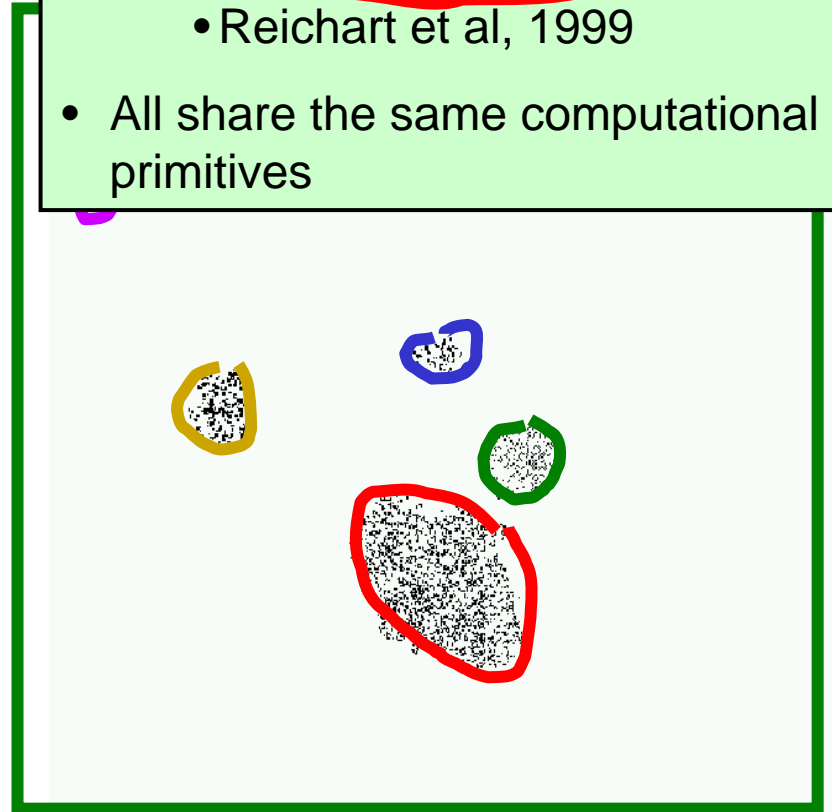
- Many possible approaches
- Examples:
 - Dasgupta and Raftery, 1998
 - Byers and Raftery, 1998
 - Cuevas et al, 2000
 - Reichart et al, 1999
- All share the same computational primitives



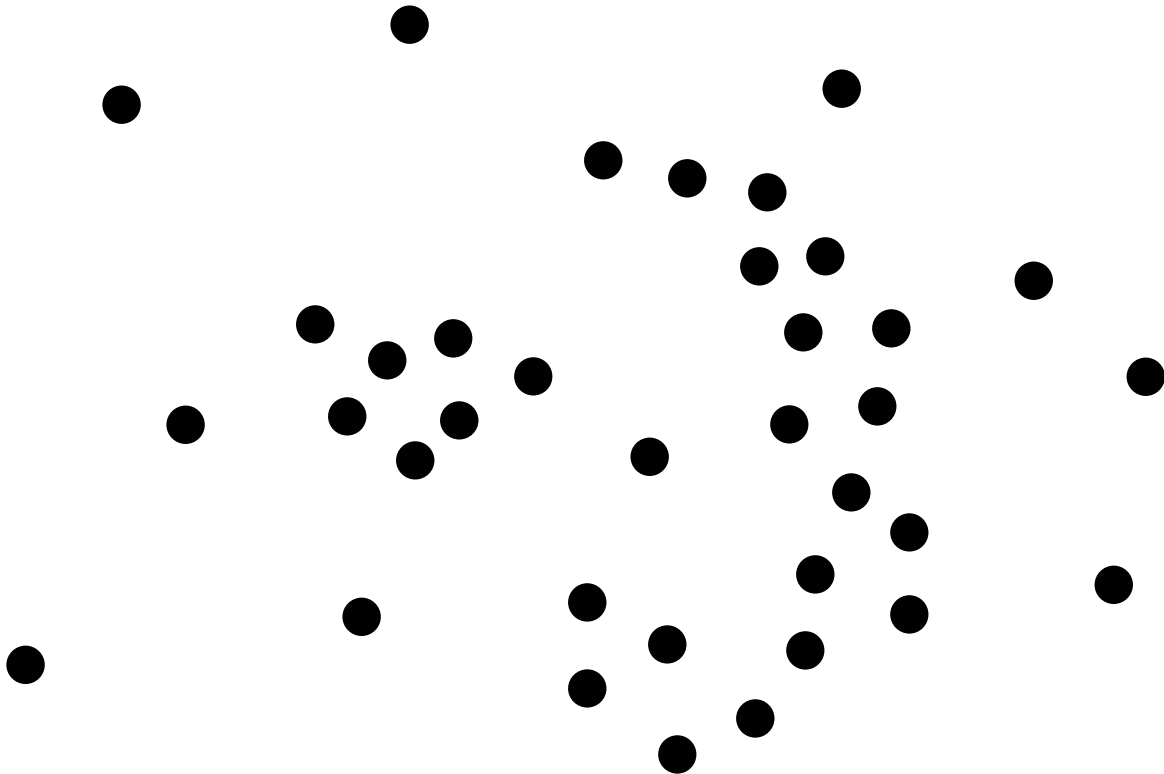
Finding the overdense regions



- Many possible approaches
- Examples:
 - Dasgupta and Raftery, 1998
 - Byers and Raftery, 1998
 - Cuevas et al, 2000
 - Reichart et al, 1999
- All share the same computational primitives

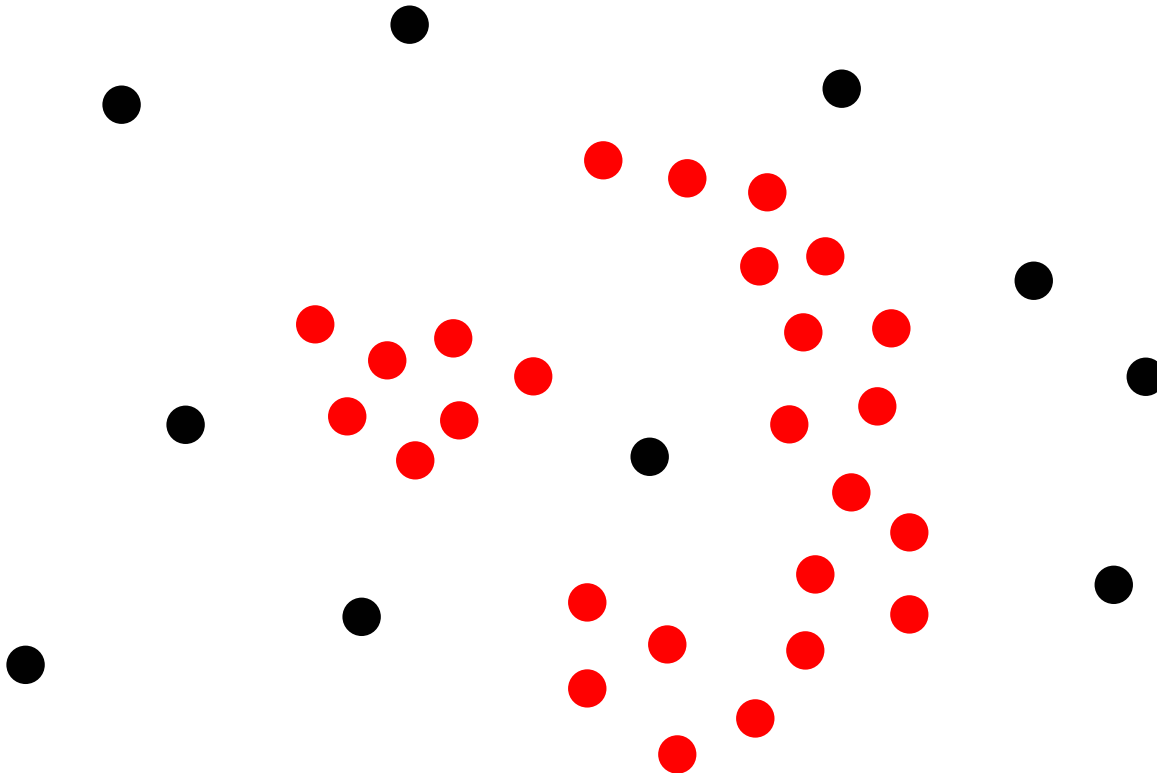


Finding the overdense regions



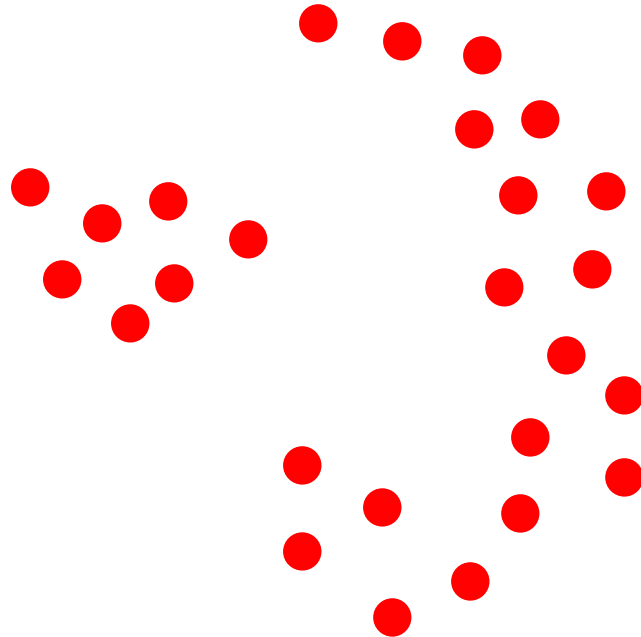
Finding the overdense regions

- Step One:
Identify the high
density points



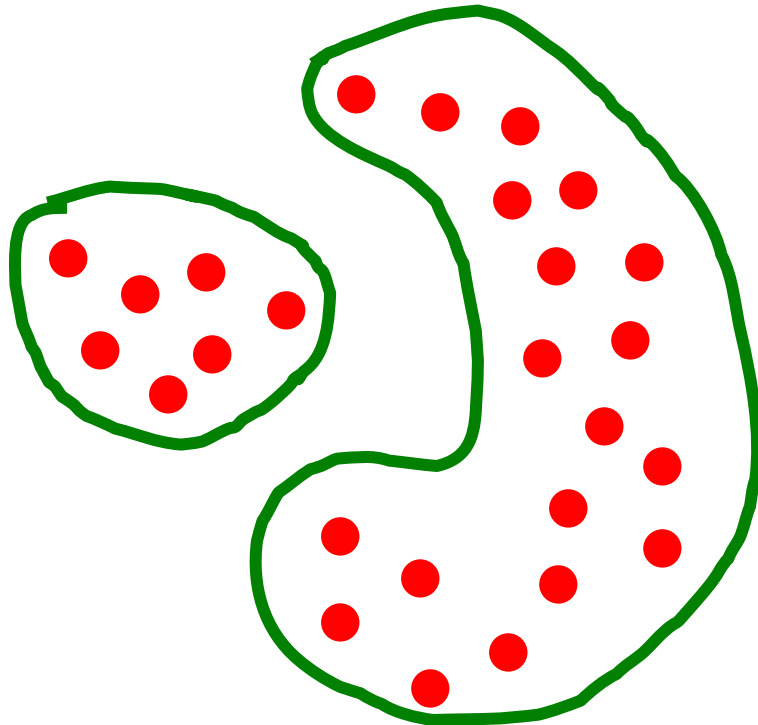
Finding the overdense regions

- Step One:
Identify the high density points
- Step Two:
Delete the rest



Finding the overdense regions

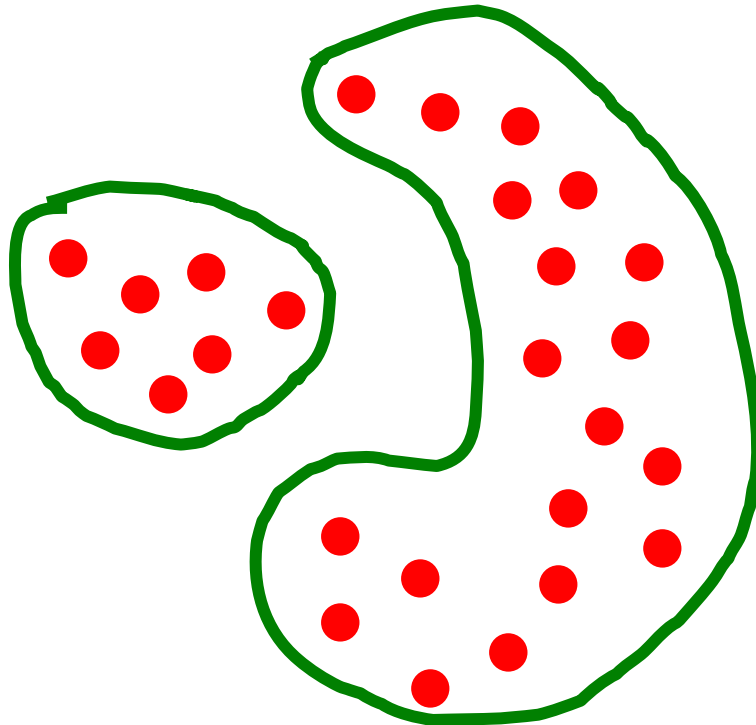
- Step One:
Identify the high density points
- Step Two:
Delete the rest
- Step Three:
Find connected components



Finding the overdense regions

CFF assumes:
Kernel Density Estimation

- Step One: Identify the high density points
- Step Two: Delete the rest
- Step Three: Find connected components

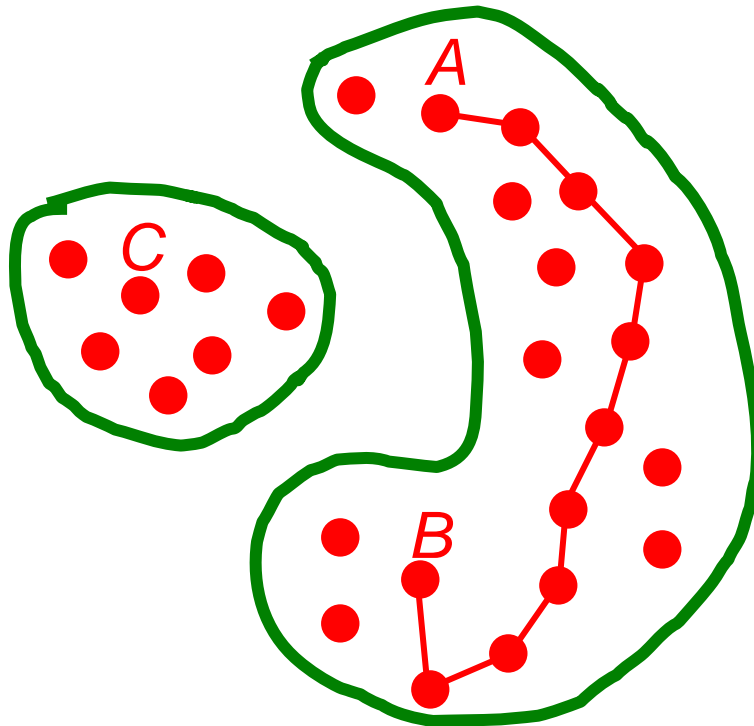


Finding the overdense regions

CFF assumes:
Kernel Density Estimation

- Step One: Identify the high density points
- Step Two: Delete the rest
- Step Three: Find connected components

CFF assumes: A and B are in same component if there's a path between them with all small steps



Finding the overdense region

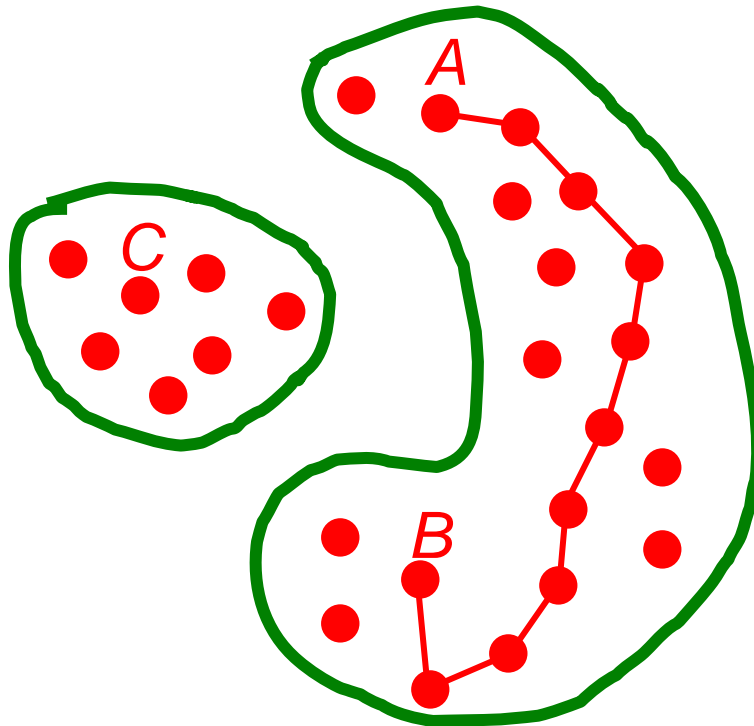
Can be done efficiently with 2-point style search

CFF assumes:
Kernel Density Estimation

- Step One: Identify the high density points
- Step Two: Delete the rest
- Step Three: Find connected components

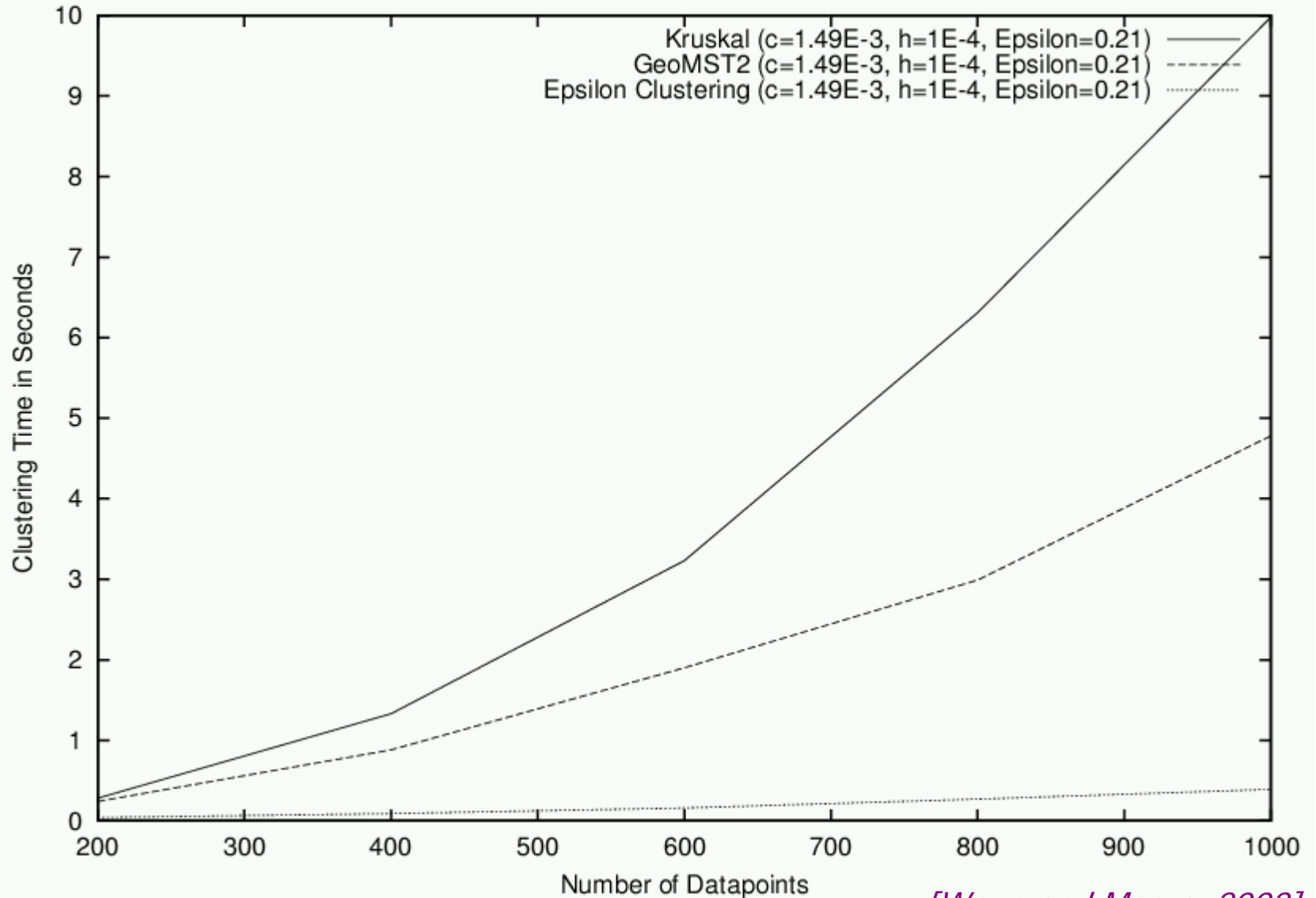
CFF assumes: A and B are in same component if there's a path between them with all small steps

Can be done efficiently with 2-point-style search plus an extra trick



Results

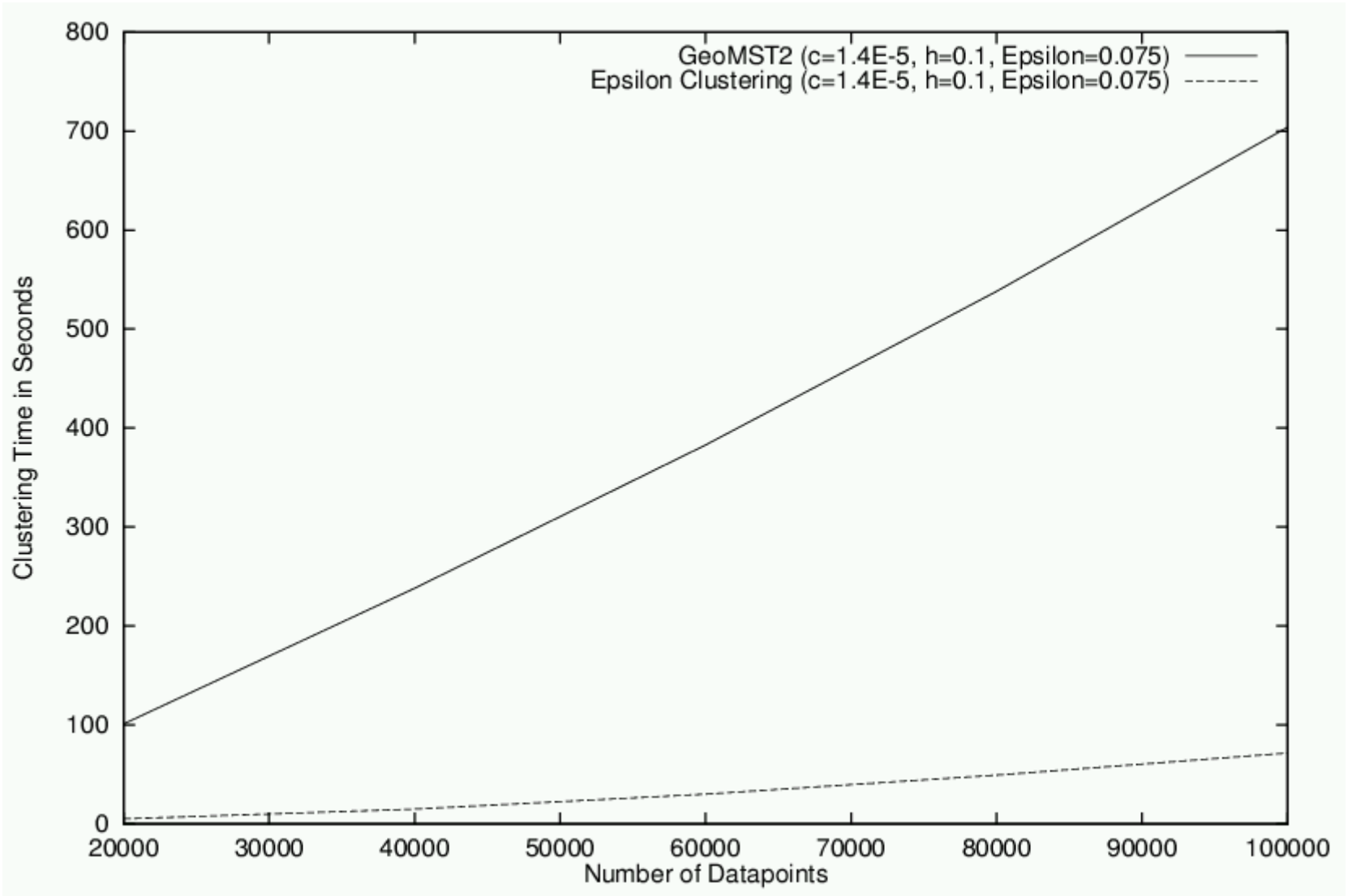
4-dimensional Sloan
Astrophysics color-space data



[Wong and Moore, 2002]

Results

4-dimensional Sloan
Astrophysics color-space data



Outline

Cached Sufficient Statistics

Ball Trees (= Metric Trees)

K-nearest neighbor with ball trees

- Very fast non-parametric classification

- skewed binary outputs

- General binary outputs

- multi-classed outputs

Very fast kernel-based statistics

- n-point computations

- clustering

- non-parametric clustering (overdensity hunting)

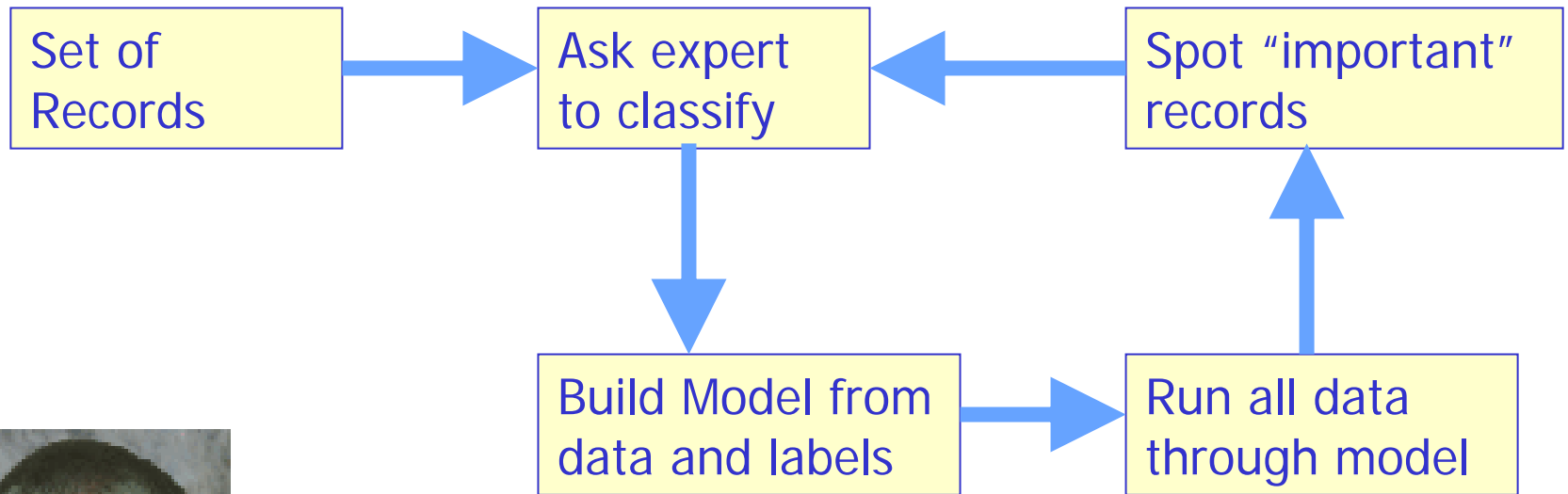


- Active learning for anomaly hunting

- GMorph: Efficient Galaxy morphology fitting

- Other Auton topics

Active Learning of Anomalies



Dan Pelleg

Anomaly GUI

File

Previous batch of images 4 - 7 Next batch of images Find Anomalies

The screenshot displays the Anomaly GUI interface. On the left, there are four vertically stacked panels, each showing a different astronomical image with a central star and a crosshair. The top panel is labeled '5' and 'N' at the top, 'E' on the left, and 'W' on the right. Below each image panel are several controls: a row of checkboxes for 'c0', 'c4', 'c5', 'c6', 'c1', 'c2', and 'c3'; a 'locked' checkbox; 'zoom in' and 'zoom out' buttons; and a 'hands off!' checkbox. To the right of these controls is a horizontal bar composed of 20 colored segments, representing the detection of anomalies. The colors of these segments vary across the four panels, indicating different detected anomalies. The top bar of the GUI is cyan, and the overall background is light gray.

c0 c4 c5 c6 c1 c2 c3

locked zoom in zoom out hands off!

c0 c4 c5 c6 c1 c2 c3

locked zoom in zoom out hands off!

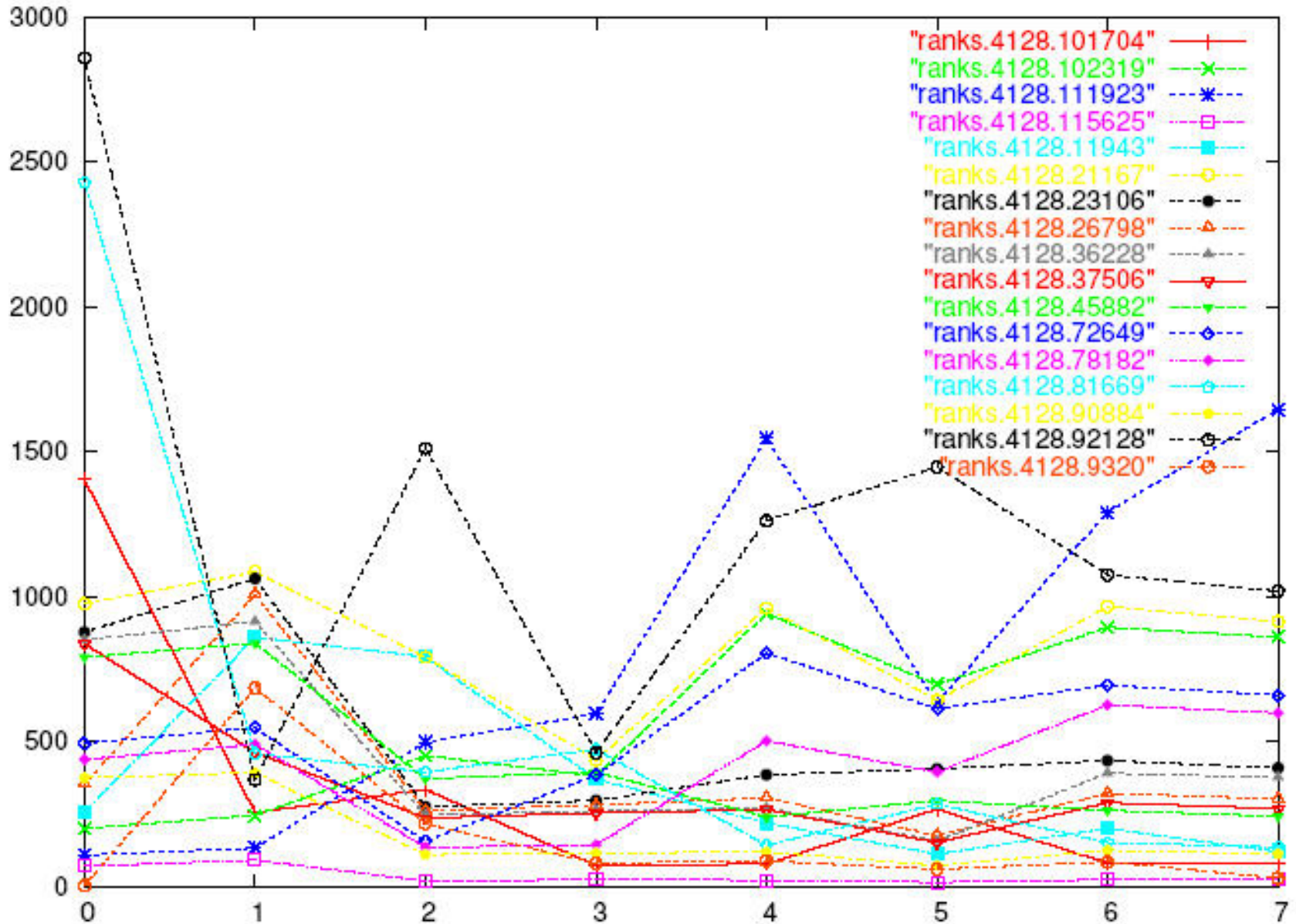
c0 c4 c5 c6 c1 c2 c3

locked zoom in zoom out hands off!

c0 c4 c5 c6 c1 c2 c3

locked zoom in zoom out hands off!

Anomaly Performance



Outline

Cached Sufficient Statistics

Ball Trees (= Metric Trees)

K-nearest neighbor with ball trees

- Very fast non-parametric classification

- skewed binary outputs

- General binary outputs

- multi-classed outputs

Very fast kernel-based statistics

- n-point computations

- clustering

- non-parametric clustering (overdensity hunting)

Active learning for anomaly hunting

▶ GMorph: Efficient Galaxy morphology fitting

Other Auton topics

GMorph: Fast Galaxy Morphology

- How do you perform 10^7 large nonlinear optimizations in practical time?
- How do you avoid local optima
- Idea: Pre-cache a “library” of solutions. Use efficient nearest neighbor to match new problems to library as seeds.
- Early tests bring galaxy morphology fits down from minutes to sub-seconds



Brigham
Anderson

Outline

Cached Sufficient Statistics

Ball Trees (= Metric Trees)

K-nearest neighbor with ball trees

- Very fast non-parametric classification

- skewed binary outputs

- General binary outputs

- multi-classed outputs

Very fast kernel-based statistics

- n-point computations

- clustering

- non-parametric clustering (overdensity hunting)

Active learning for anomaly hunting

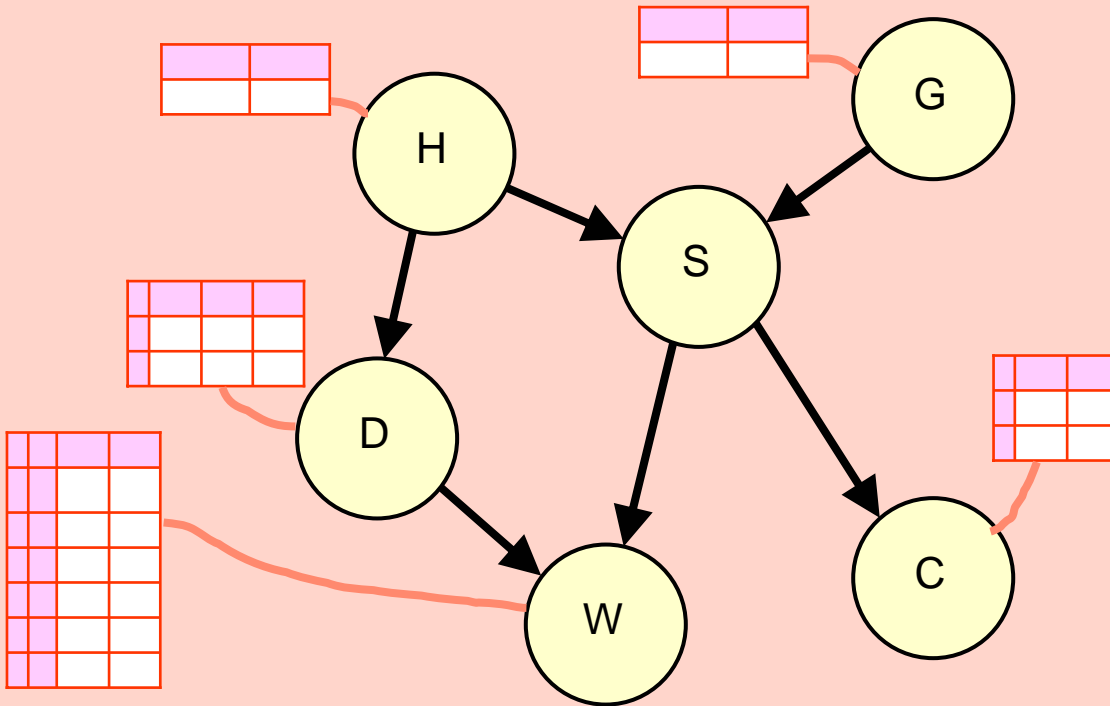
GMorph: Efficient Galaxy morphology fitting



Other Auton topics

Other Relevant Auton Topics

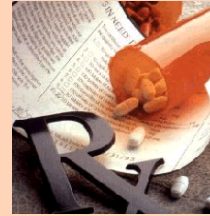
Bayesian Networks



Other Relevant Auton Topics

Bayesian Networks

"What's strange about recent events?"



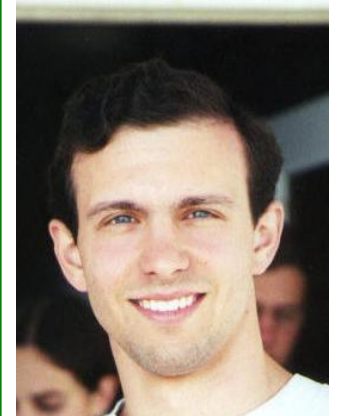
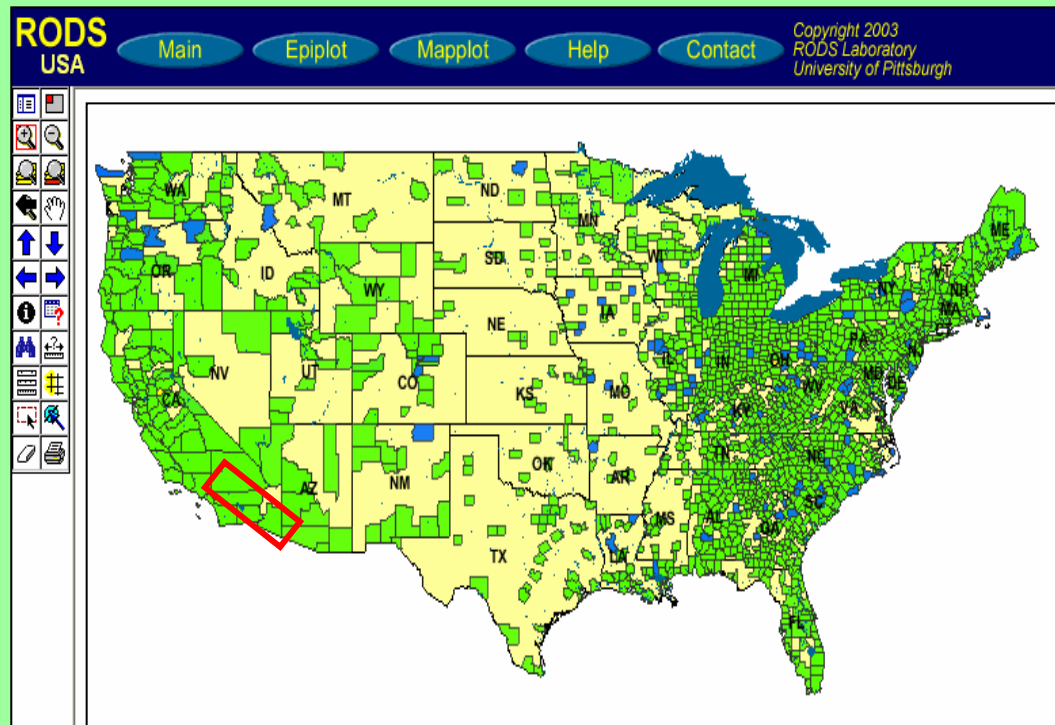
Weng-Keen
Wong

Other Relevant Auton Topics

Bayesian Networks

“What’s strange about recent events?”

Spatial Scan Statistics



Daniel Neill

[Neill, Moore and Wagner, 2004]

Other Relevant Auton Topics

Bayesian Networks

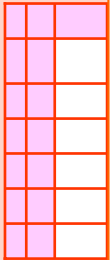
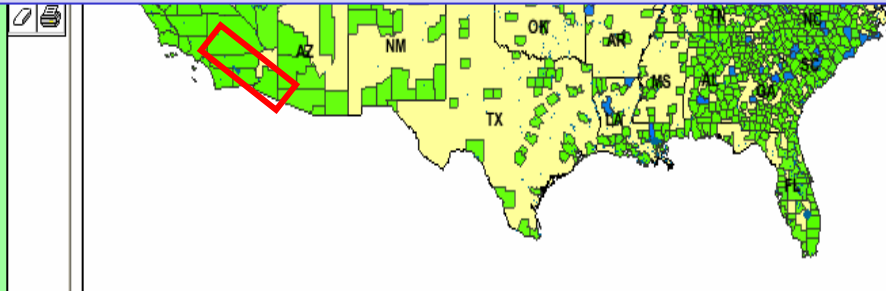
"What's strange about recent events?"

Spatial Scan Statistics

Massively multiple target tracking



Jeremy Kubica



Conclusions

- Geometry can help tractability of Massive Statistical Data Analysis
- Cached sufficient statistics are one approach
- Papers, tutorials, software, examples:

www.autonlab.org