

# Next-Generation Capabilities for Large-Scale Scientific Visualization

SIAM Conference on Parallel Processing for Scientific Computing  
February 15, 2012

Kenneth Moreland, Nathan Fabian  
Sandia National Laboratories

Berk Geveci, Utkarsh Ayachit  
Kitware Inc.

James Ahrens  
Los Alamos National Laboratory

\*\*\*\* Release Markings \*\*\*\*

# Collaborators

- **Sandia National Laboratories**

- Kenneth Moreland
- Nathan Fabian
- David Thompson
- Ron Oldfield

- **Los Alamos National Laboratories**

- James Ahrens
- Jonathan Woodring

- **Oak Ridge National Laboratory**

- Scott Klasky
- Norbert Podhorszki

- **Argonne National Laboratory**

- Venkatram Vishwanath
- Mark Hereld
- Michael E. Papka

- **Kitware, Inc.**

- Berk Geveci
- Utkarsh Ayachit
- Andrew C. Bauer
- Pat Marion
- Sebastien Jourdain
- David DeMarle

- **University of Colorado at Boulder**

- Michel Rasquin
- Kenneth E. Jansen

- **Rutgers University**

- Ciprian Docan
- Manish Parashar

# Outline (Hide)

- Overview history: Onyx to Exascale
- ParaView
  - General capabilities
  - Remote visualization
- In situ
  - Workflow, In transit
- Exascale challenges
  - What changes
  - Dax

# 1990's

Before consumer market impacted graphics. Single memory/multipipe machines (SGI)

# 2000's

- 2005: Specialized Vis cluster - Distributed memory, commodity clusters tightly coupled with compute platform.
  - Specialized HW (graphics cards, I/O, memory)
- 2010: Distributed memory capability clusters – 'Running on the platform'
  - Highly constrained memory, no specialized HW

# 2020's

- Constrained by power
  - Need: Coupled analysis
- Exascale breaks everything
  - New programming models
  - Billion-way parallelism



## RedRoSE/BlackRoSE

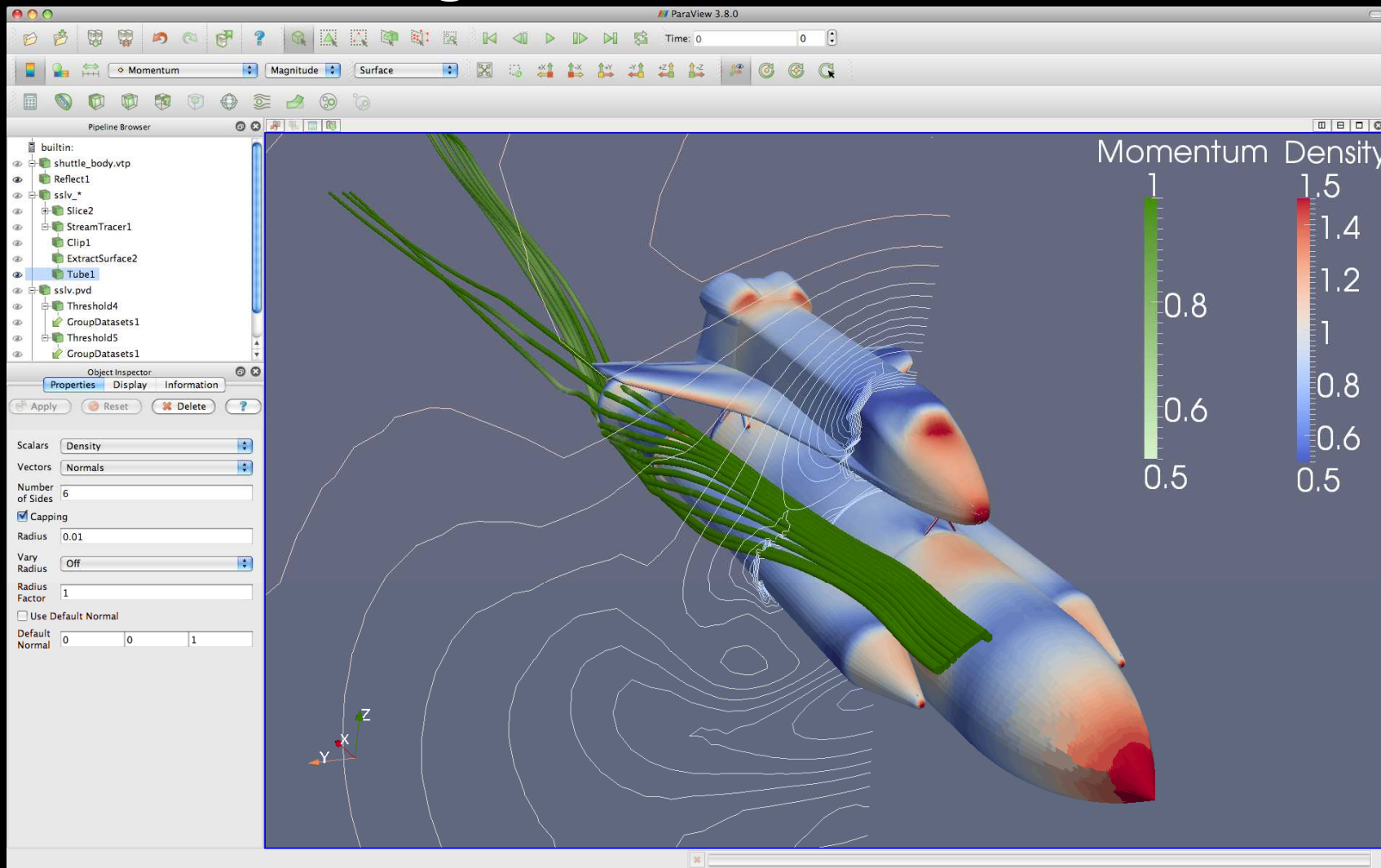


Cielo

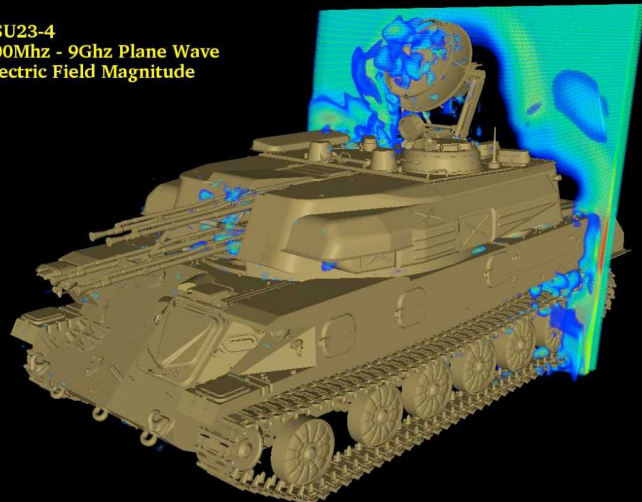
## Exascale Platform



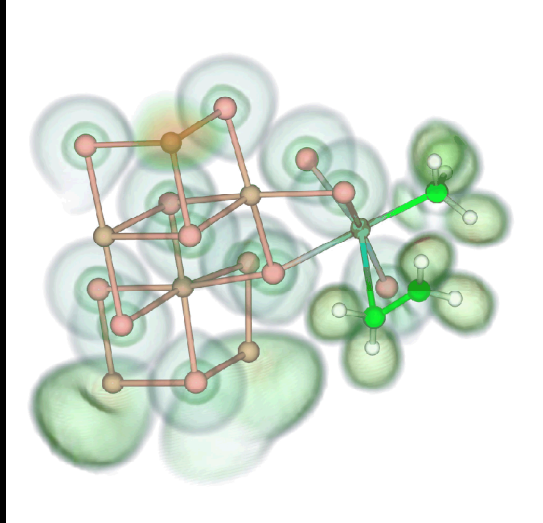
# ParaView: an End User Tool for Large-Scale Visualization



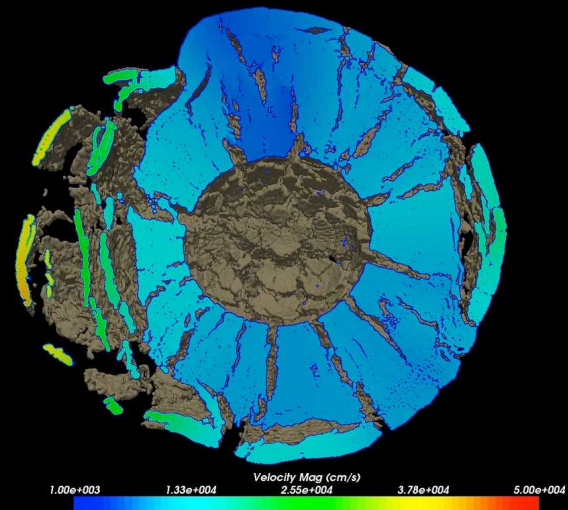
ZSU23-4  
100Mhz - 9Ghz Plane Wave  
Electric Field Magnitude



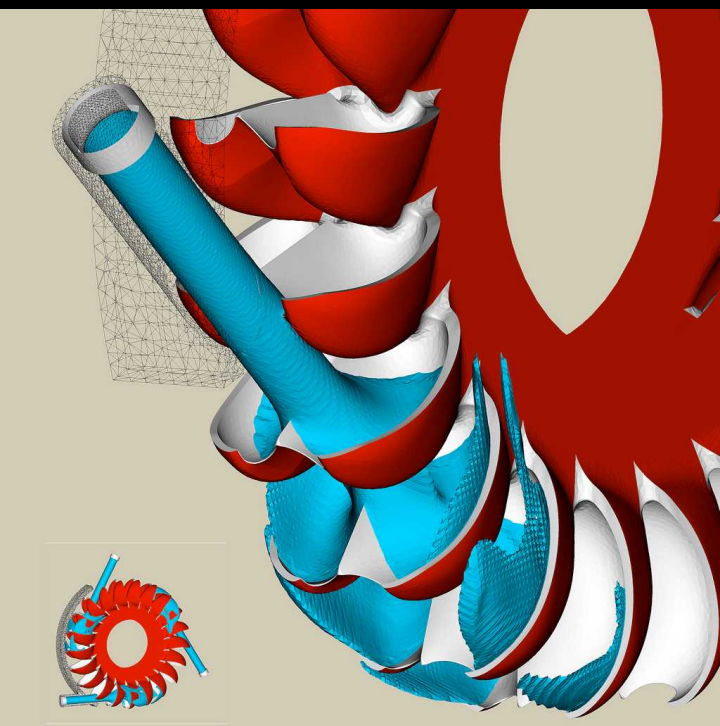
Jerry Clarke, US Army Research Laboratory



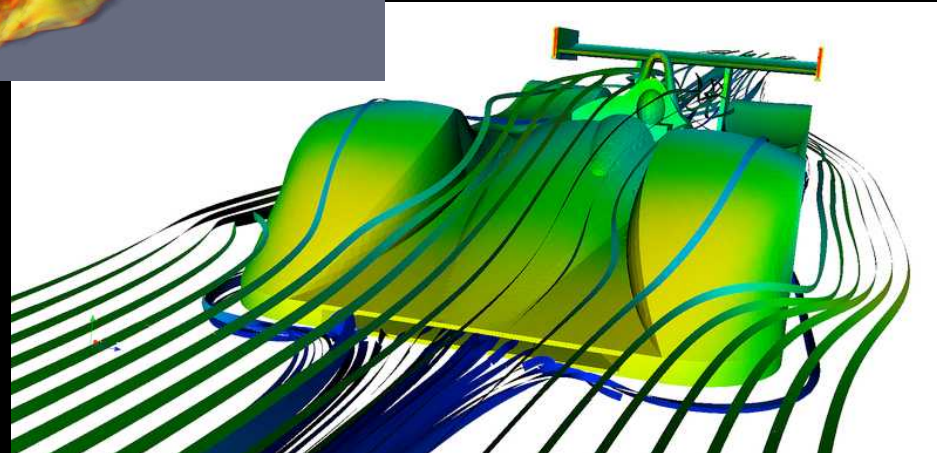
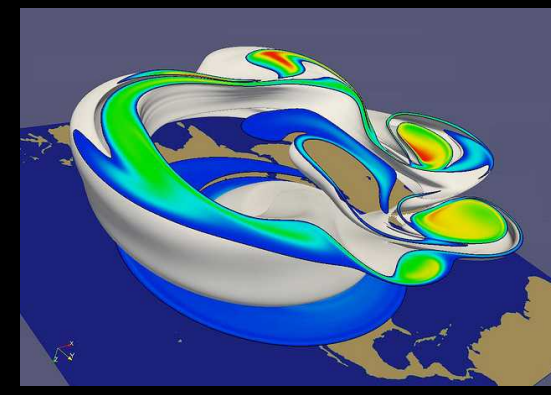
Swiss National Supercomputing Centre



Velocity Mag (cm/s)  
1.00e+003 1.33e+004 2.55e+004 3.78e+004 5.00e+004



Swiss National Supercomputing Centre



Renato N. Elias, NACAD/COPPE/UFRJ, Rio de Janeiro, Brazil

# 1990's

Before consumer market impacted graphics. Single memory/multipipe machines (SGI)

# 2000's

- 2005: Specialized Vis cluster - Distributed memory, commodity clusters tightly coupled with compute platform.
  - Specialized HW (graphics cards, I/O, memory)
- 2010: Distributed memory capability clusters – 'Running on the platform'
  - Highly constrained memory, no specialized HW

# 2020's

- Constrained by power
  - Need: Coupled analysis
- Exascale breaks everything
  - New programming models
  - Billion-way parallelism



## RedRoSE/BlackRoSE



Cielo

## Exascale Platform



# Slide of Doom

	2010	“2018”	Factor Change
System Peak	2 Pf/s	1 Ef/s	500
Power	6 MW	20 MW	3
System Memory	0.3 PB	10 PB	33
Node Performance	0.125 Gf/s	10 Tf/s	80
Node Memory BW	25 GB/s	400 GB/s	16
Node Concurrency	12 cpus	1,000 cpus	83
Interconnect BW	1.5 GB/s	50 GB/s	33
System Size (nodes)	20 K nodes	1 M nodes	50
Total Concurrency	225 K	1 B	4,444
Storage	15 PB	300 PB	20
Input/Output bandwidth	0.2 TB/s	20 TB/s	100

# Slide of Doom

	2010	"2018"	Factor Change
System Peak	2 Pf/s	1 Ef/s	
Power	6 MW	20 MW	3
System Memory	0.3 PB	10 PB	33
Node Performance	0.125 Gf/s	10 Tf/s	80
Node Memory BW	25 GB/s	400 GB/s	16
Node Concurrency	12 cpus	1,000 cpus	83
Interconnect BW	1.5 GB/s	50 GB/s	33
System Size (nodes)	20 K nodes	1 M nodes	50
Total Concurrency	225 K	1 B	4,444
Storage	15 PB	300 PB	20
Input/Output bandwidth	0.2 TB/s	20 TB/s	

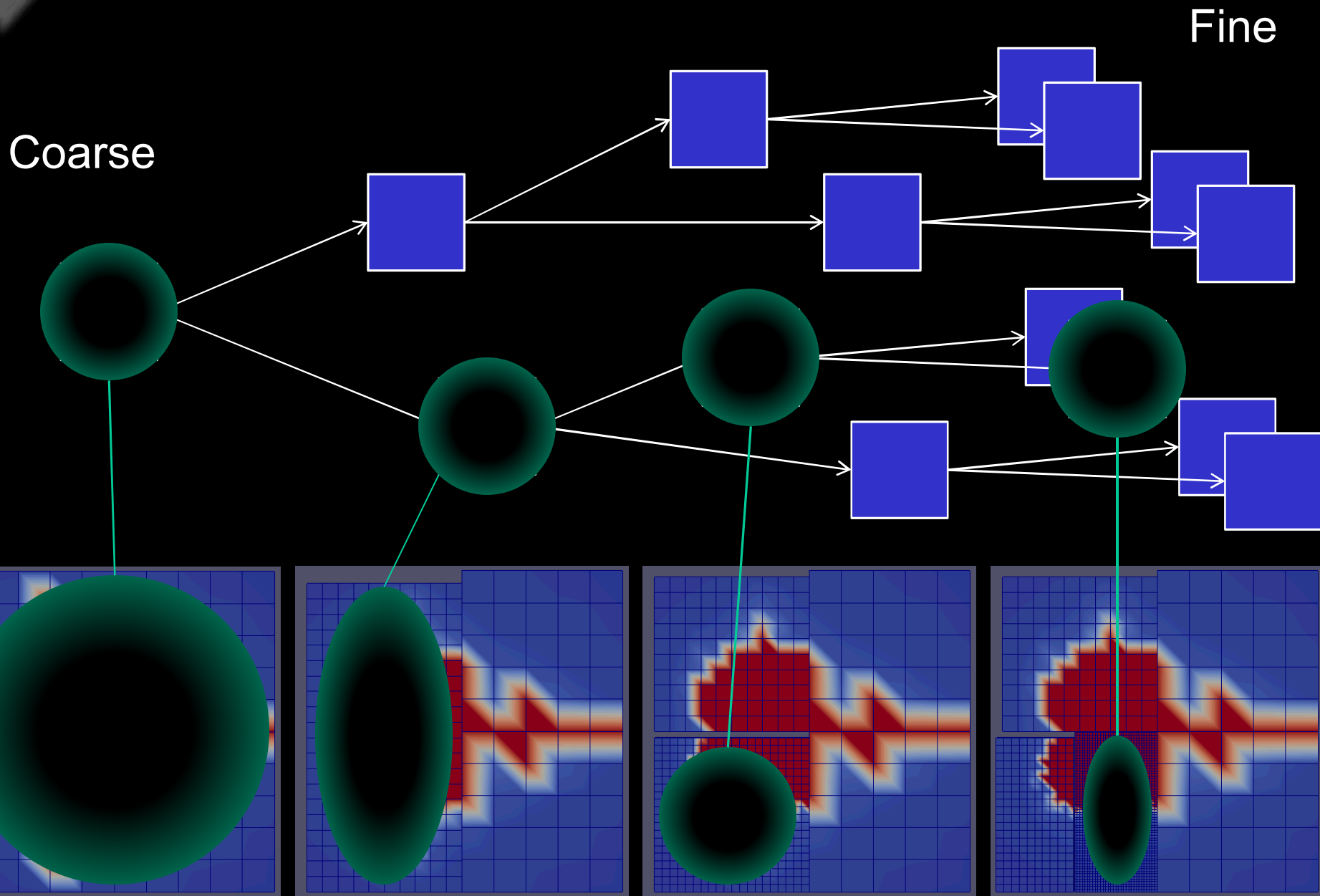
# Extreme scale computing

- Trends
  - More FLOPS
  - More concurrency
  - Comparatively less storage, I/O bandwidth
- ASCI purple (49 TB/140 GB/s)
- JaguarPF (300 TB/200 GB/s)
  - Most people get < 5 GB/sec at scale



Systems	2009	2011	2015	2018
System Peak Flops/s	2 Peta	20 Peta	100-200 Peta	1 Exa
System Memory	0.3 PB	1 PB	5 PB	10 PB
Node Performance	125 GF	200 GF	400 GF	1-10 TF
Node Memory BW	25 GB/s	40 GB/s	100 GB/s	200-400 GB/s
Node Concurrency	12	32	0(100)	0(1000)
Interconnect BW	1.5 GB/s	10 GB/s	25 GB/s	50 GB/s
System Size (Nodes)	18,700	100,000	500,000	0(Million)
Total Concurrency	225,000	3 Million	50 Million	0(Billion)
Storage	15 PB	30 PB	150 PB	300 PB
I/O	0.2 TB/s	2 TB/s	10 TB/s	20 TB/s
MTTI	Days	Days	Days	0(1Day)
Power	6 MW	~10 MW	~10 MW	~20 MW

# Prioritized Adaptive Streaming



PointCenteredData

Surface With Edge

Pipeline Browser

- builtin:
- CFC11.t.t0.1\_42l\_nccs01.005905.raw

Object Inspector

Properties Display Information

Apply Reset Delete ?

Spacing 1 1 1

Whole Extent	0	3599
	0	2399
	0	41

Swap Endian

Refinement Inspector

Per Object controls

Restart

Lock Refinement

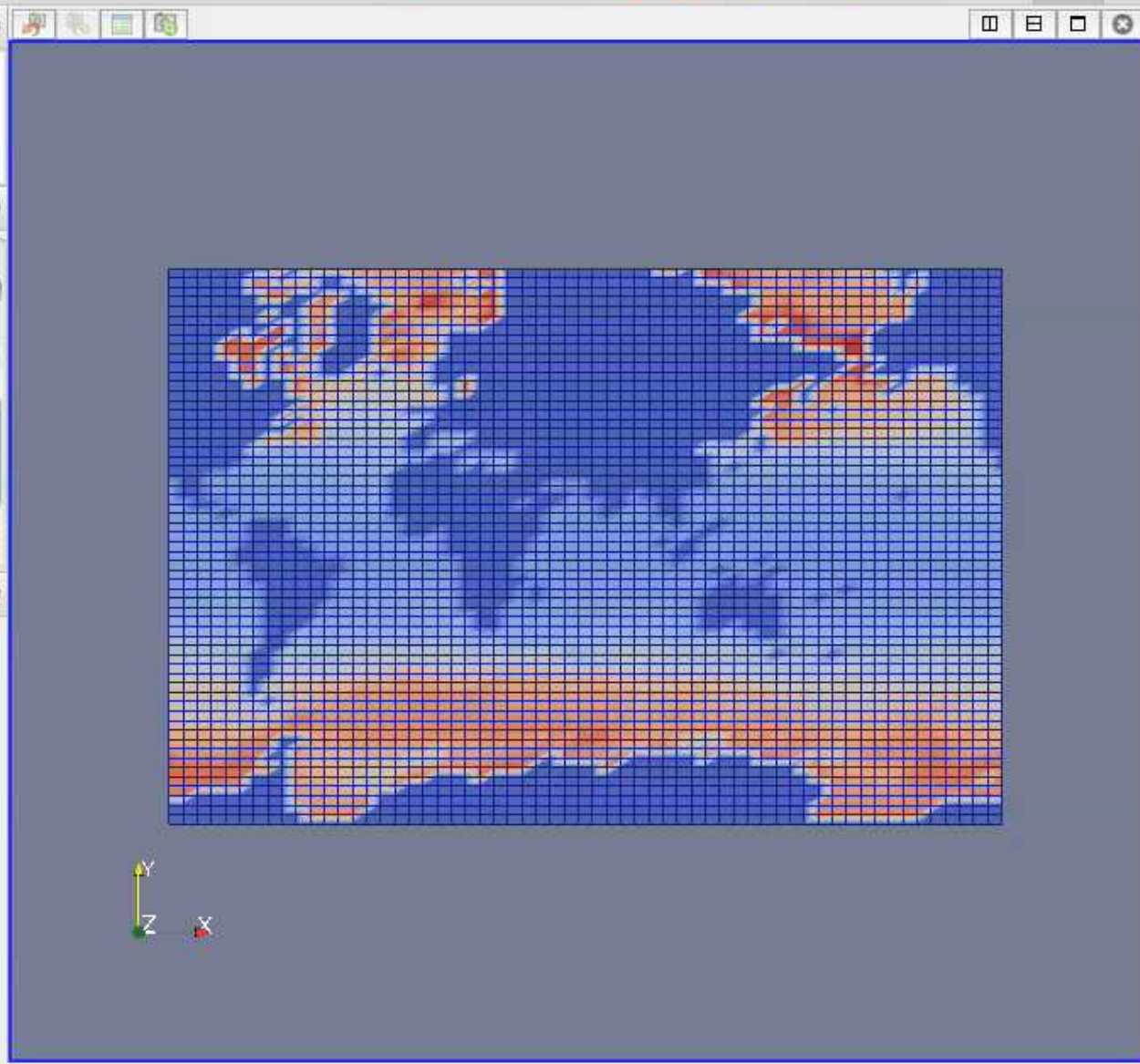
Show Piece Bounds

Depth cut off -1

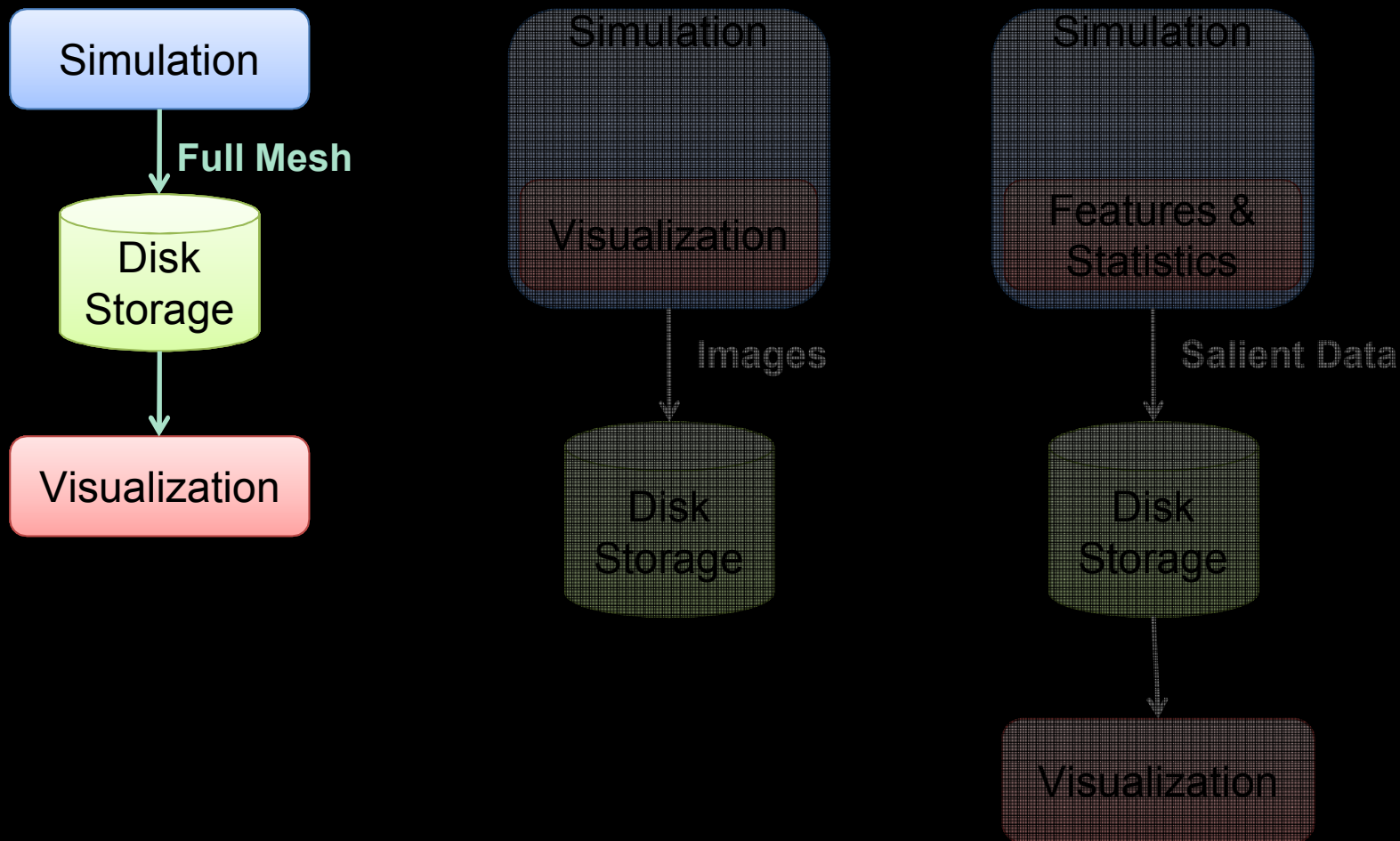
Global controls

INTERRUPT

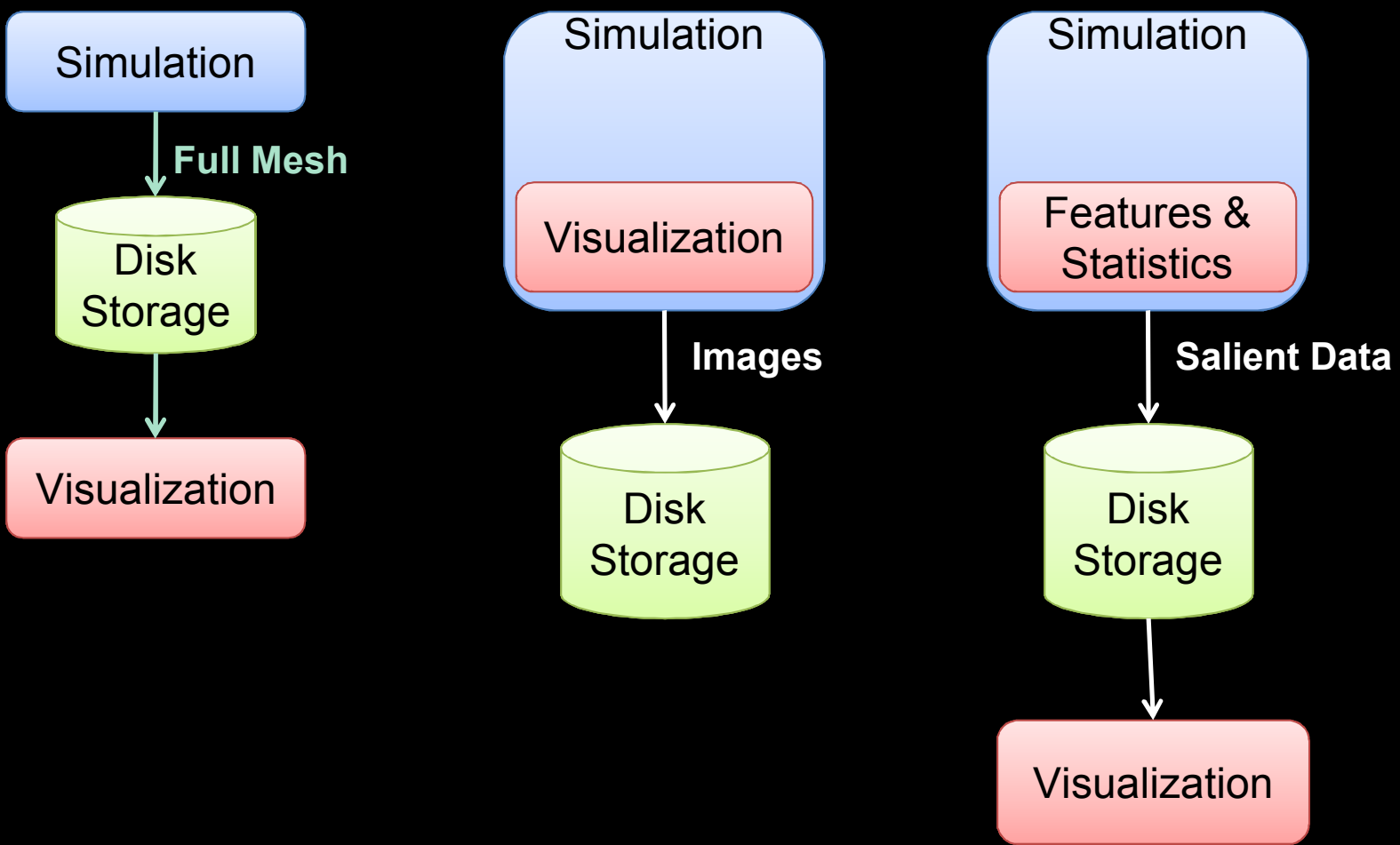
MANUAL - +

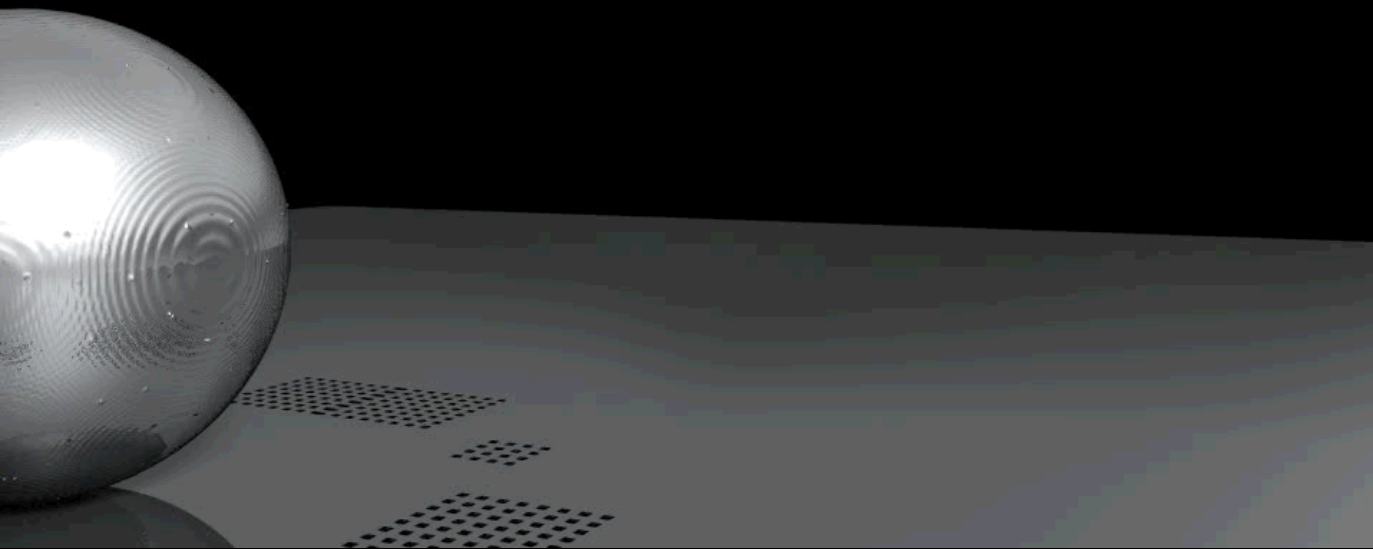


# Diskless Visualization

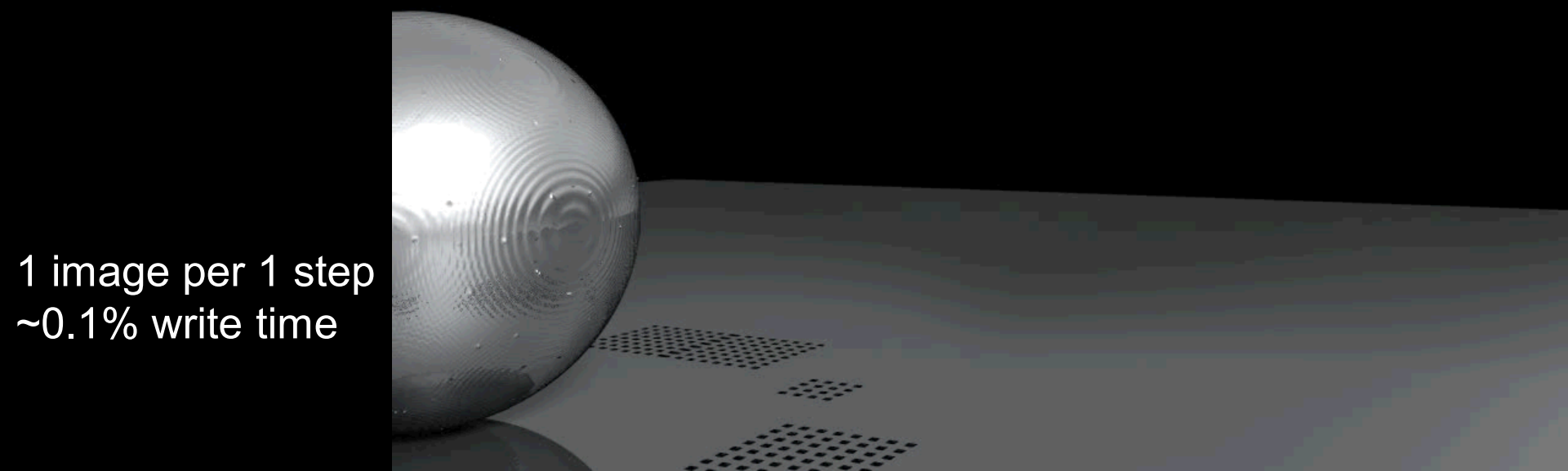


# Diskless Visualization





1 mesh per 100 steps  
~1% write time



1 image per 1 step  
~0.1% write time

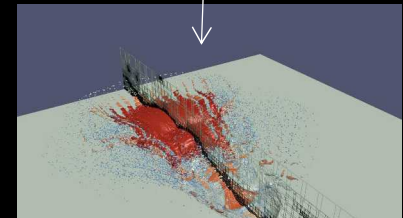
Simulation

ParaView  
Coproprocessing

Simulation

ParaView  
Coproprocessing

Output  
Processed  
Data



Rendered Images

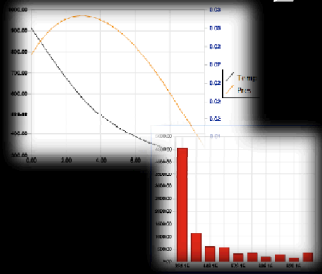
Simulation

ParaView  
Coproprocessing

Output  
Processed  
Data

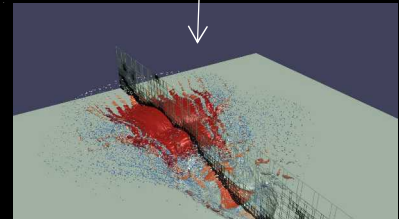
Time	Step	Value 1	Value 2	Value 3	Value 4	Value 5
0.000000	1	0.000000	0.000000	0.000000	0.000000	0.000000
0.000000	2	0.000000	0.000000	0.000000	0.000000	0.000000
0.000000	3	0.000000	0.000000	0.000000	0.000000	0.000000
0.000000	4	0.000000	0.000000	0.000000	0.000000	0.000000
0.000000	5	0.000000	0.000000	0.000000	0.000000	0.000000
0.000000	6	0.000000	0.000000	0.000000	0.000000	0.000000
0.000000	7	0.000000	0.000000	0.000000	0.000000	0.000000
0.000000	8	0.000000	0.000000	0.000000	0.000000	0.000000
0.000000	9	0.000000	0.000000	0.000000	0.000000	0.000000
0.000000	10	0.000000	0.000000	0.000000	0.000000	0.000000

Statistics



Line Series

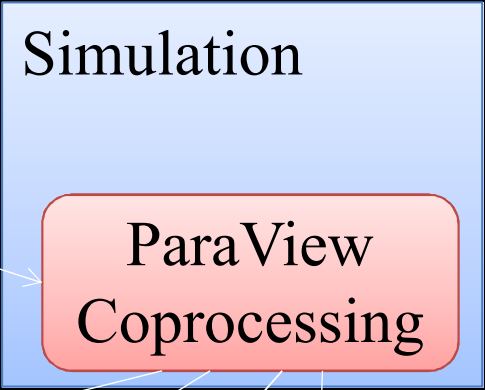
Polygonal Surfaces  
Field Data



Rendered Images

```
# Create the reader and set the filename.
reader = servermanager.sources.Reader(FileNames=path)
view = servermanager.CreateRenderView()
repr = servermanager.CreateRepresentation(reader, view)
reader.UpdatePipeline()
dataInfo = reader.GetDataInformation()
pDInfo = dataInfo.GetPointDataInformation()
arrayInfo = pDInfo.GetArrayInformation("displacement9")
if arrayInfo:
    # get the range for the magnitude of displacement9
    range = arrayInfo.GetComponentRange(-1)
    lut = servermanager.rendering.PVLookupTable()
    lut.RGBPoints = [range[0], 0.0, 0.0, 1.0,
                    range[1], 1.0, 0.0, 0.0]
    lut.VectorMode = "Magnitude"
    repr.LookupTable = lut
    repr.ColorArrayName = "displacement9"
    repr.ColorAttributeType = "POINT_DATA"
```

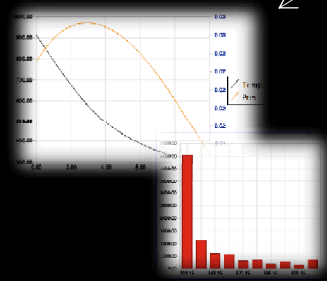
Augmented  
script in input  
deck.



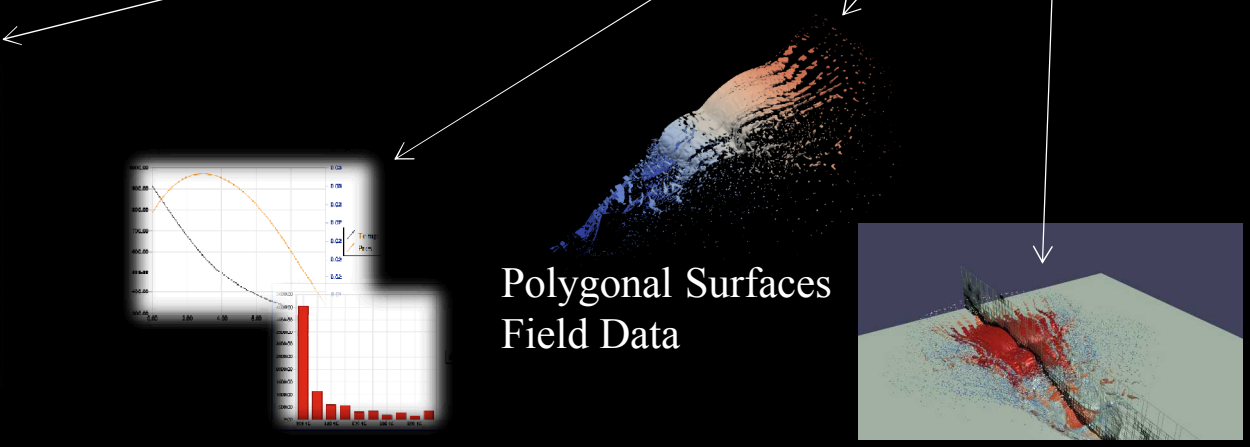
Output  
Processed  
Data

A large table with multiple columns and rows of numerical data, representing statistical information extracted from the simulation.

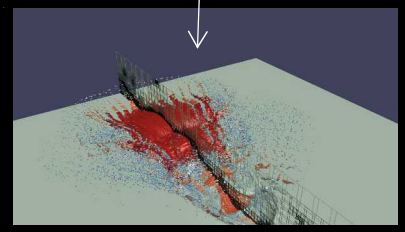
Statistics



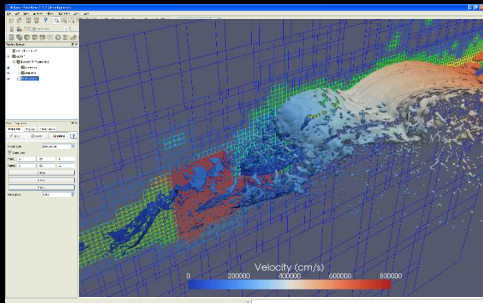
Line Series



Polygonal Surfaces  
Field Data



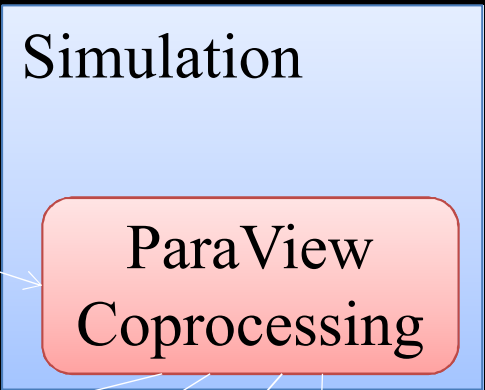
Rendered Images



Script Export

```
# Create the reader and set the filename.
reader = servermanager.sources.Reader(FileNames=path)
view = servermanager.CreateRenderView()
repr = servermanager.CreateRepresentation(reader, view)
reader.UpdatePipeline()
dataInfo = reader.GetDataInformation()
pDInfo = dataInfo.GetPointDataInformation()
arrayInfo = pDInfo.GetArrayInformation("displacement9")
if arrayInfo:
    # get the range for the magnitude of displacement9
    range = arrayInfo.GetComponentRange(-1)
    lut = servermanager.rendering.PVLookupTable()
    lut.RGBPoints = [range[0], 0.0, 0.0, 1.0,
                    range[1], 1.0, 0.0, 0.0]
    lut.VectorMode = "Magnitude"
    repr.LookupTable = lut
    repr.ColorArrayName = "displacement9"
    repr.ColorAttributeType = "POINT_DATA"
```

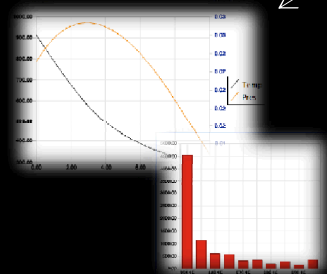
Augmented  
script in input  
deck.



Output  
Processed  
Data

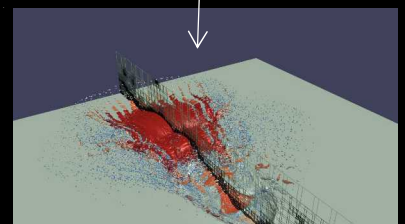
Statistics

Time	Step	Time	Step	Time	Step	Time	Step
0.000000	0	0.000000	0	0.000000	0	0.000000	0
0.000000	1	0.000000	1	0.000000	1	0.000000	1
0.000000	2	0.000000	2	0.000000	2	0.000000	2
0.000000	3	0.000000	3	0.000000	3	0.000000	3
0.000000	4	0.000000	4	0.000000	4	0.000000	4
0.000000	5	0.000000	5	0.000000	5	0.000000	5
0.000000	6	0.000000	6	0.000000	6	0.000000	6
0.000000	7	0.000000	7	0.000000	7	0.000000	7
0.000000	8	0.000000	8	0.000000	8	0.000000	8
0.000000	9	0.000000	9	0.000000	9	0.000000	9
0.000000	10	0.000000	10	0.000000	10	0.000000	10
0.000000	11	0.000000	11	0.000000	11	0.000000	11
0.000000	12	0.000000	12	0.000000	12	0.000000	12
0.000000	13	0.000000	13	0.000000	13	0.000000	13
0.000000	14	0.000000	14	0.000000	14	0.000000	14
0.000000	15	0.000000	15	0.000000	15	0.000000	15
0.000000	16	0.000000	16	0.000000	16	0.000000	16
0.000000	17	0.000000	17	0.000000	17	0.000000	17
0.000000	18	0.000000	18	0.000000	18	0.000000	18
0.000000	19	0.000000	19	0.000000	19	0.000000	19
0.000000	20	0.000000	20	0.000000	20	0.000000	20
0.000000	21	0.000000	21	0.000000	21	0.000000	21
0.000000	22	0.000000	22	0.000000	22	0.000000	22
0.000000	23	0.000000	23	0.000000	23	0.000000	23
0.000000	24	0.000000	24	0.000000	24	0.000000	24
0.000000	25	0.000000	25	0.000000	25	0.000000	25
0.000000	26	0.000000	26	0.000000	26	0.000000	26
0.000000	27	0.000000	27	0.000000	27	0.000000	27
0.000000	28	0.000000	28	0.000000	28	0.000000	28
0.000000	29	0.000000	29	0.000000	29	0.000000	29
0.000000	30	0.000000	30	0.000000	30	0.000000	30
0.000000	31	0.000000	31	0.000000	31	0.000000	31
0.000000	32	0.000000	32	0.000000	32	0.000000	32
0.000000	33	0.000000	33	0.000000	33	0.000000	33
0.000000	34	0.000000	34	0.000000	34	0.000000	34
0.000000	35	0.000000	35	0.000000	35	0.000000	35
0.000000	36	0.000000	36	0.000000	36	0.000000	36
0.000000	37	0.000000	37	0.000000	37	0.000000	37
0.000000	38	0.000000	38	0.000000	38	0.000000	38
0.000000	39	0.000000	39	0.000000	39	0.000000	39
0.000000	40	0.000000	40	0.000000	40	0.000000	40
0.000000	41	0.000000	41	0.000000	41	0.000000	41
0.000000	42	0.000000	42	0.000000	42	0.000000	42
0.000000	43	0.000000	43	0.000000	43	0.000000	43
0.000000	44	0.000000	44	0.000000	44	0.000000	44
0.000000	45	0.000000	45	0.000000	45	0.000000	45
0.000000	46	0.000000	46	0.000000	46	0.000000	46
0.000000	47	0.000000	47	0.000000	47	0.000000	47
0.000000	48	0.000000	48	0.000000	48	0.000000	48
0.000000	49	0.000000	49	0.000000	49	0.000000	49
0.000000	50	0.000000	50	0.000000	50	0.000000	50
0.000000	51	0.000000	51	0.000000	51	0.000000	51
0.000000	52	0.000000	52	0.000000	52	0.000000	52
0.000000	53	0.000000	53	0.000000	53	0.000000	53
0.000000	54	0.000000	54	0.000000	54	0.000000	54
0.000000	55	0.000000	55	0.000000	55	0.000000	55
0.000000	56	0.000000	56	0.000000	56	0.000000	56
0.000000	57	0.000000	57	0.000000	57	0.000000	57
0.000000	58	0.000000	58	0.000000	58	0.000000	58
0.000000	59	0.000000	59	0.000000	59	0.000000	59
0.000000	60	0.000000	60	0.000000	60	0.000000	60
0.000000	61	0.000000	61	0.000000	61	0.000000	61
0.000000	62	0.000000	62	0.000000	62	0.000000	62
0.000000	63	0.000000	63	0.000000	63	0.000000	63
0.000000	64	0.000000	64	0.000000	64	0.000000	64
0.000000	65	0.000000	65	0.000000	65	0.000000	65
0.000000	66	0.000000	66	0.000000	66	0.000000	66
0.000000	67	0.000000	67	0.000000	67	0.000000	67
0.000000	68	0.000000	68	0.000000	68	0.000000	68
0.000000	69	0.000000	69	0.000000	69	0.000000	69
0.000000	70	0.000000	70	0.000000	70	0.000000	70
0.000000	71	0.000000	71	0.000000	71	0.000000	71
0.000000	72	0.000000	72	0.000000	72	0.000000	72
0.000000	73	0.000000	73	0.000000	73	0.000000	73
0.000000	74	0.000000	74	0.000000	74	0.000000	74
0.000000	75	0.000000	75	0.000000	75	0.000000	75
0.000000	76	0.000000	76	0.000000	76	0.000000	76
0.000000	77	0.000000	77	0.000000	77	0.000000	77
0.000000	78	0.000000	78	0.000000	78	0.000000	78
0.000000	79	0.000000	79	0.000000	79	0.000000	79
0.000000	80	0.000000	80	0.000000	80	0.000000	80
0.000000	81	0.000000	81	0.000000	81	0.000000	81
0.000000	82	0.000000	82	0.000000	82	0.000000	82
0.000000	83	0.000000	83	0.000000	83	0.000000	83
0.000000	84	0.000000	84	0.000000	84	0.000000	84
0.000000	85	0.000000	85	0.000000	85	0.000000	85
0.000000	86	0.000000	86	0.000000	86	0.000000	86
0.000000	87	0.000000	87	0.000000	87	0.000000	87
0.000000	88	0.000000	88	0.000000	88	0.000000	88
0.000000	89	0.000000	89	0.000000	89	0.000000	89
0.000000	90	0.000000	90	0.000000	90	0.000000	90
0.000000	91	0.000000	91	0.000000	91	0.000000	91
0.000000	92	0.000000	92	0.000000	92	0.000000	92
0.000000	93	0.000000	93	0.000000	93	0.000000	93
0.000000	94	0.000000	94	0.000000	94	0.000000	94
0.000000	95	0.000000	95	0.000000	95	0.000000	95
0.000000	96	0.000000	96	0.000000	96	0.000000	96
0.000000	97	0.000000	97	0.000000	97	0.000000	97
0.000000	98	0.000000	98	0.000000	98	0.000000	98
0.000000	99	0.000000	99	0.000000	99	0.000000	99
0.000000	100	0.000000	100	0.000000	100	0.000000	100

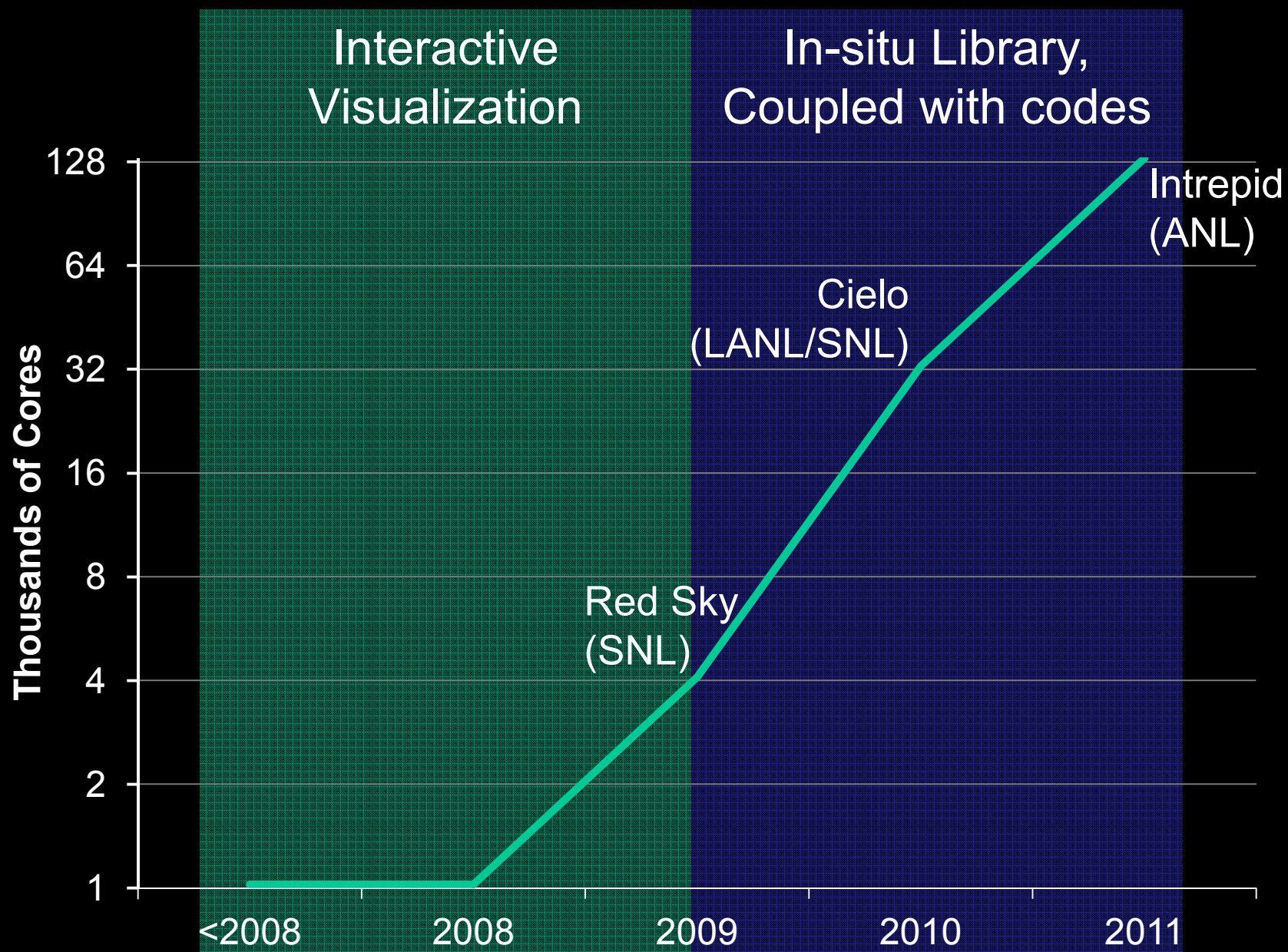


Line Series

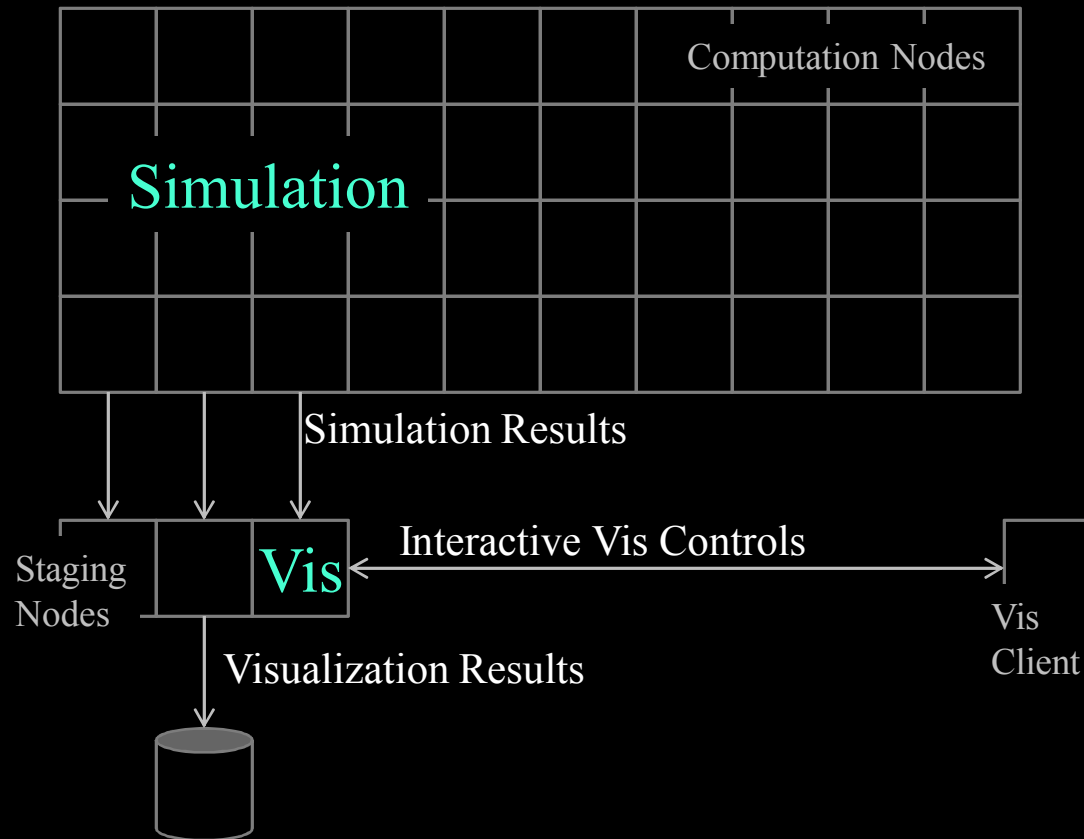
Polygonal Surfaces  
Field Data



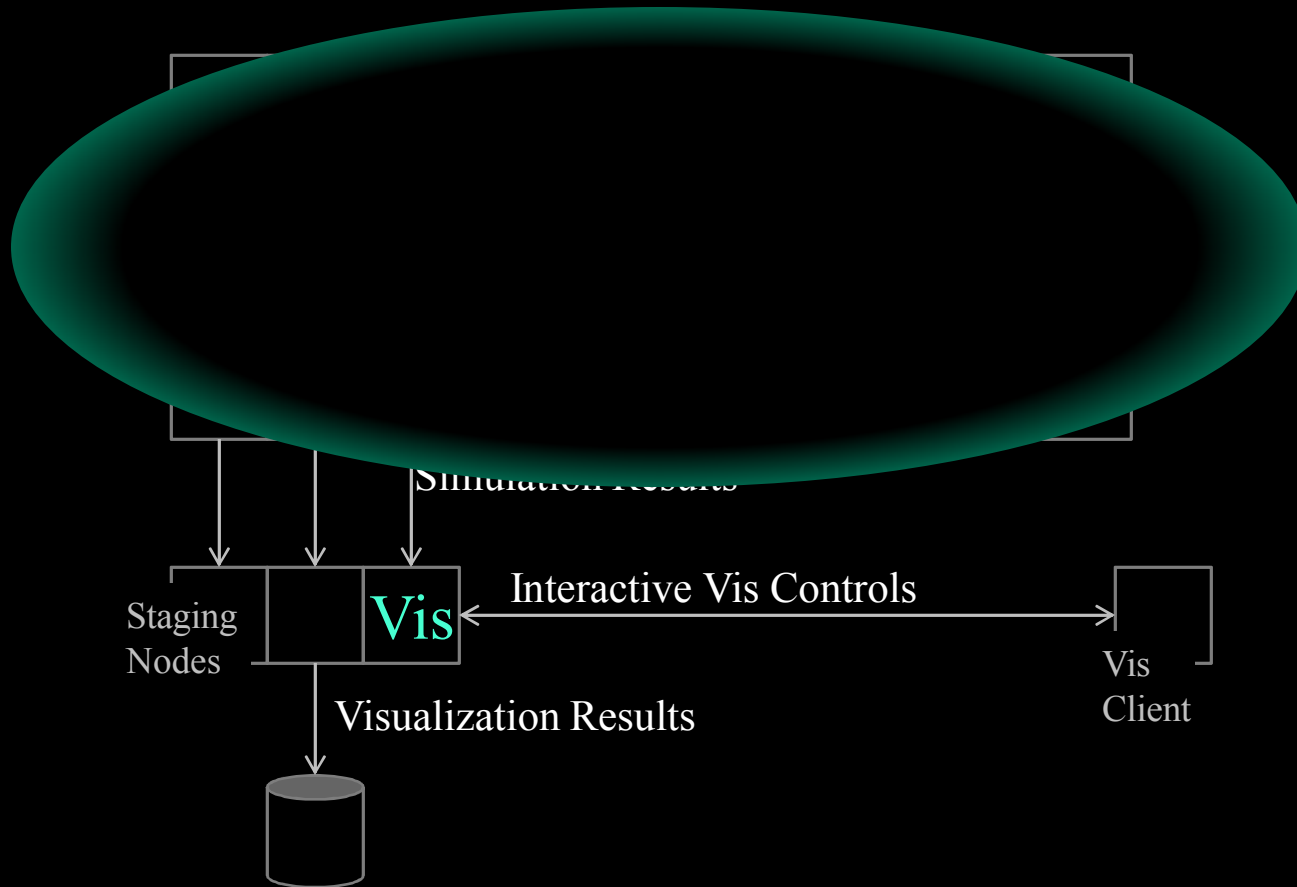
Rendered Images



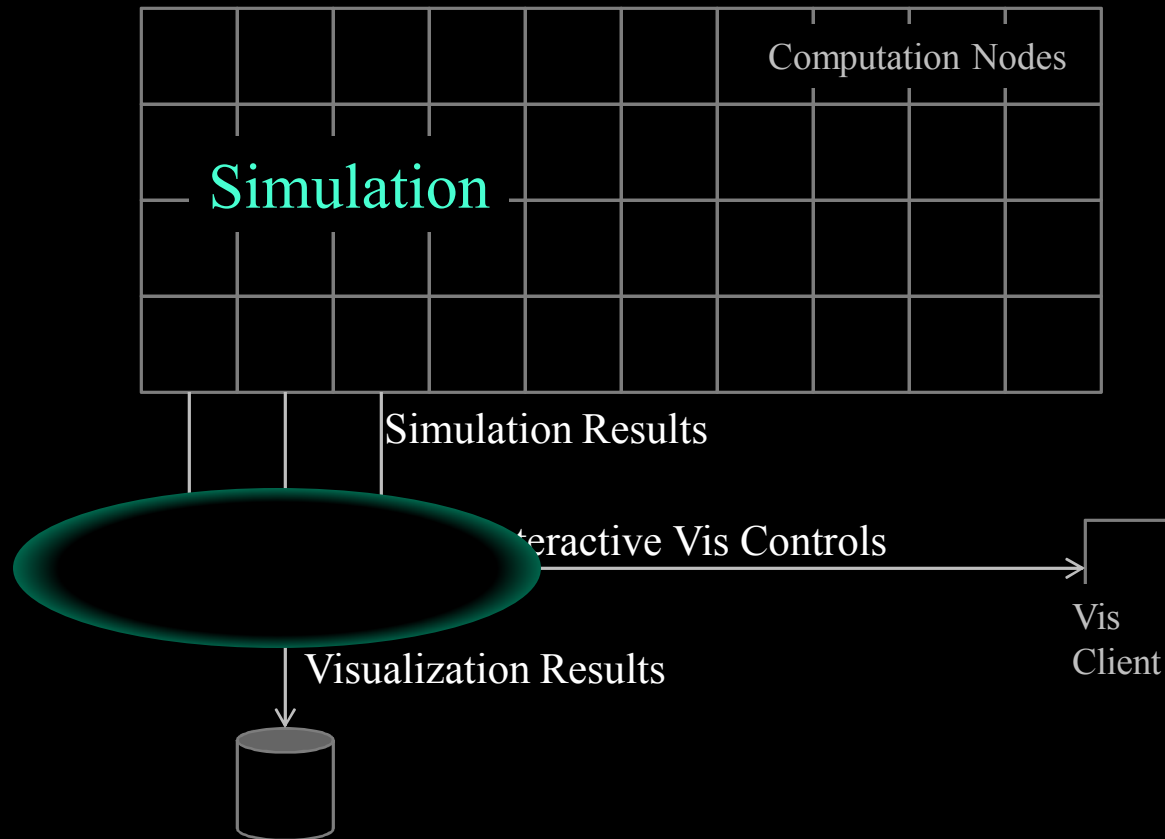
# In Transit Visualization



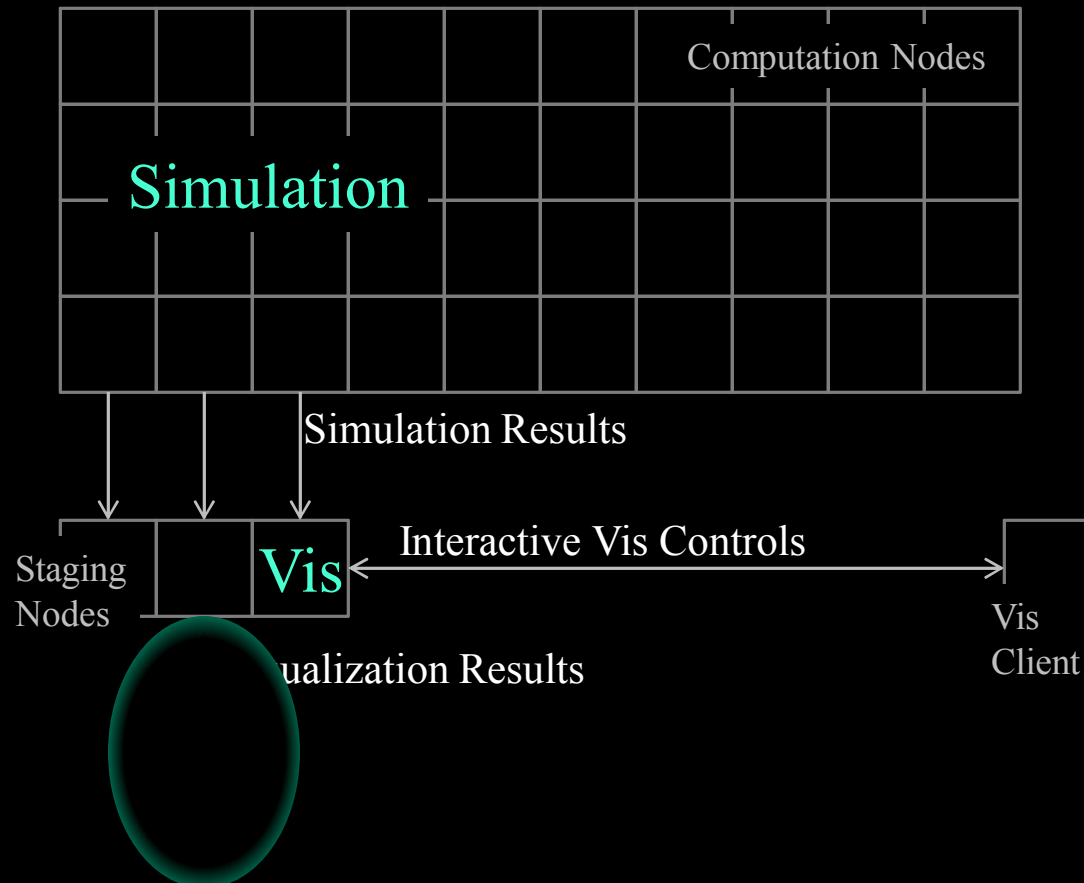
# In Transit Visualization



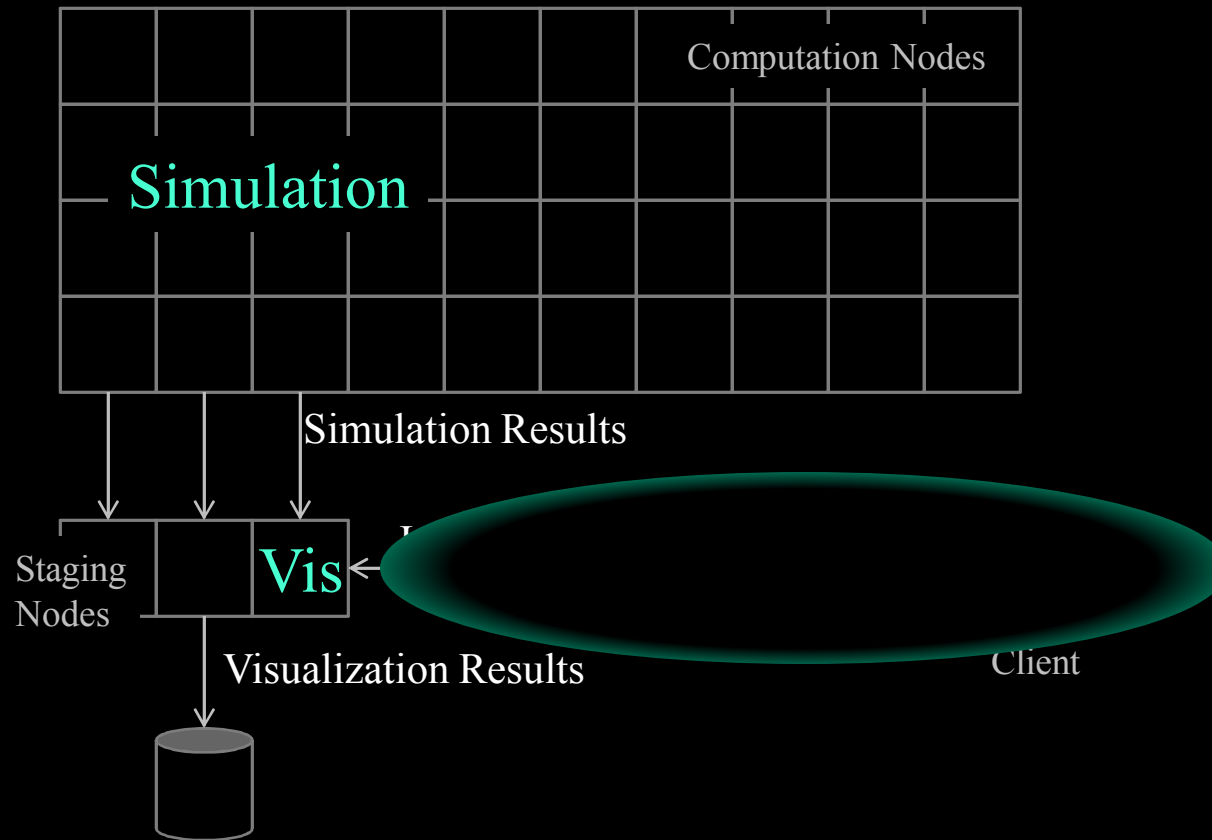
# In Transit Visualization

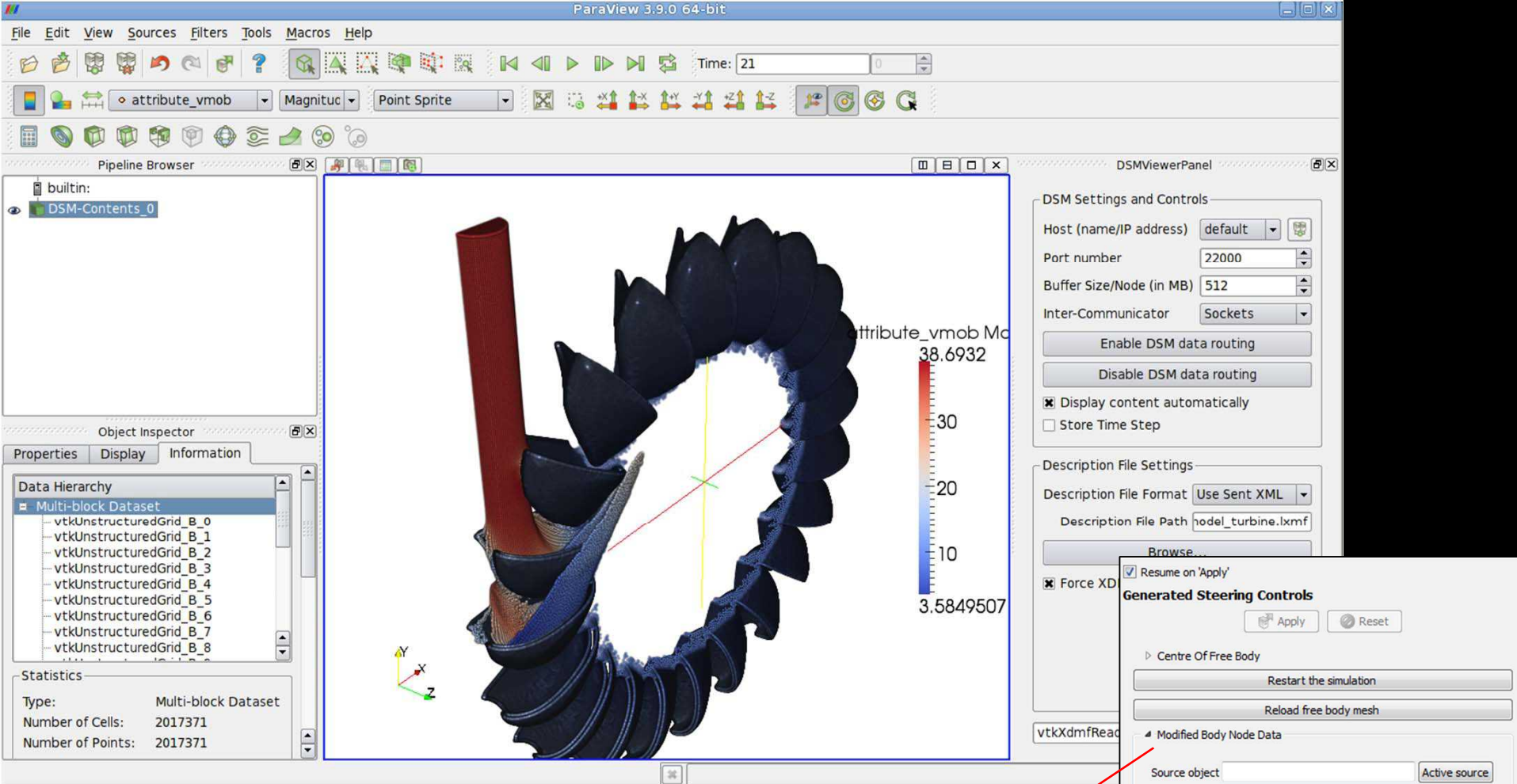


# In Transit Visualization



# In Transit Visualization

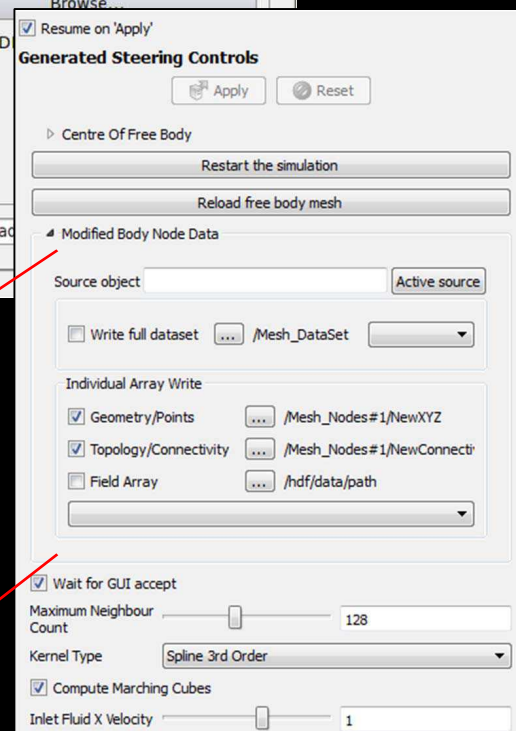


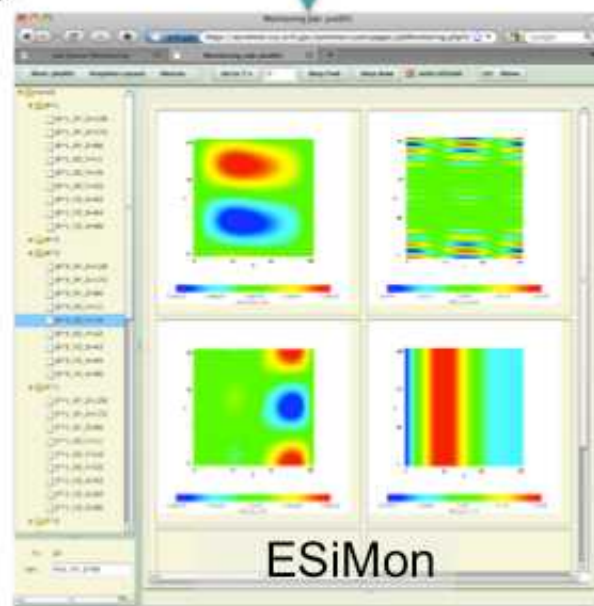
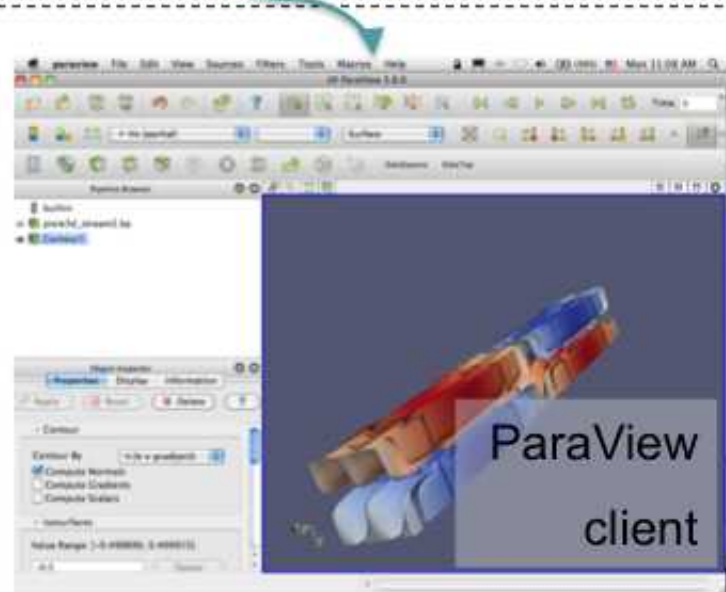
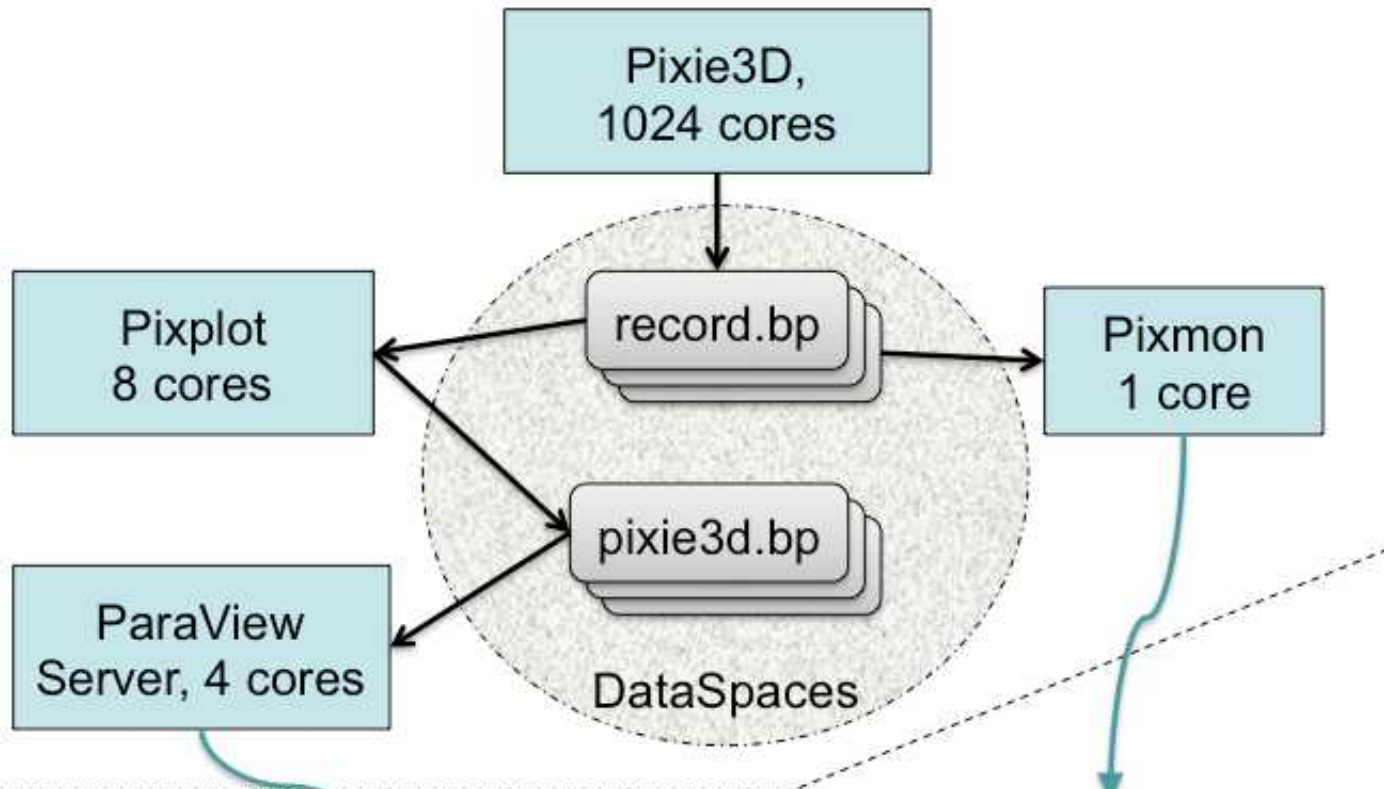


```

<DataExportProperty name="ModifiedBodyNodes"
  command="SetSteeringArray"
  label="Modified Body Node Data">
  <DataExportDomain name="data_export"
    full_path="/Mesh_DataSet"
    geometry_path="/Mesh_Nodes#1/NewXYZ"
    topology_path="/Mesh_Nodes#1/NewCo..."
    command_property="ReloadFreeBodyMesh" />
</DataExportProperty>

```





# Slide of Doom

	2010	"2018"	Factor Change
System Peak	2 Pf/s	1 Ef/s	500
Power	6 MW	20 MW	3
System Memory	0.3 PB	10 PB	33
Node Performance	0.125 Gf/s	10 Tf/s	80
Node Memory BW	25 GB/s	400 GB/s	16
Node Concurrency	12 cpus	1,000 cpus	83
Interconnect BW	1.5 GB/s	50 GB/s	33
System Size (nodes)	20 K nodes	1 M nodes	50
Total Concurrency	225 K	1 B	4,444
Storage	15 PB	300 PB	20
Input/Output bandwidth	0.2 TB/s	20 TB/s	100

# Slide of Doom

	2010	"2018"	Factor Change
System Peak	2 Pf/s	1 Ef/s	500
Power	6 MW	20 MW	3
System Memory	0.3 PB	10 PB	33
Node Performance	0.125 Gf/s	10 Tf/s	80
Node Memory BW	25 GB/s	400 GB/s	16
Node Concurrency	12 cpus	1,000 cpus	83
Interconnect BW	1.5 GB/s	50 GB/s	33
System Size (nodes)	20 K nodes	1 M nodes	50
Total Concurrency	225 K	1 B	444
Storage	15 PB	300 PB	20
Input/Output bandwidth	0.2 TB/s	20 TB/s	100

# Naïve Scaling

	2010	“2018”	Factor Change
System Memory	0.3 PB	10 PB	33
Total Concurrency	225 K	1 B	4,444

## MPI Only?

Vis object code + state: 20MB

On Jaguar: 20MB × 200,000 processes = 4TB

On Exascale: 20MB × 1 billion processes = 20PB !

# Naïve Scaling

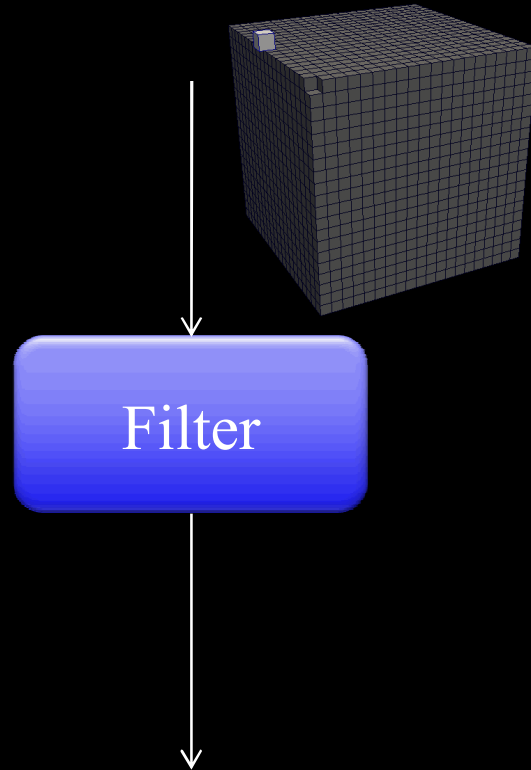
	2010	“2018”	Factor Change
System Memory	0.3 PB	10 PB	33
Total Concurrency	225 K	1 B	4,444

Visualization pipeline too heavyweight?

On Jaguar: 1 trillion cells → 5 million cells/thread

On Exascale: 30 trillion cells → 30K cells/thread

# Revisiting the Pipeline



- Lightweight Object
- Serial Execution
- No explicit partitioning
- No access to larger structures
- No state
- Changing programming models

# Dax Makes it Easy

```
int vtkCellDerivatives::RequestData(...)
{
    ...[allocate output arrays]...
    ...[validate inputs]...
    for (cellId=0; cellId < numCells; cellId++)
    {
        ...[update progress]...
        input->GetCell(cellId, cell);
        subId
            = cell->GetParametricCenter(pcoords);
        inScalars->GetTuples(cell->PointIds,
                            cellScalars);
        scalars = cellScalars->GetPointer(0);
        cell->Derivatives(subId,
                        pcoords,
                        scalars,
                        1,
                        derivs);
        outGradients->SetTuple(cellId, derivs);
    }
    ...[cleanup]...
}
```

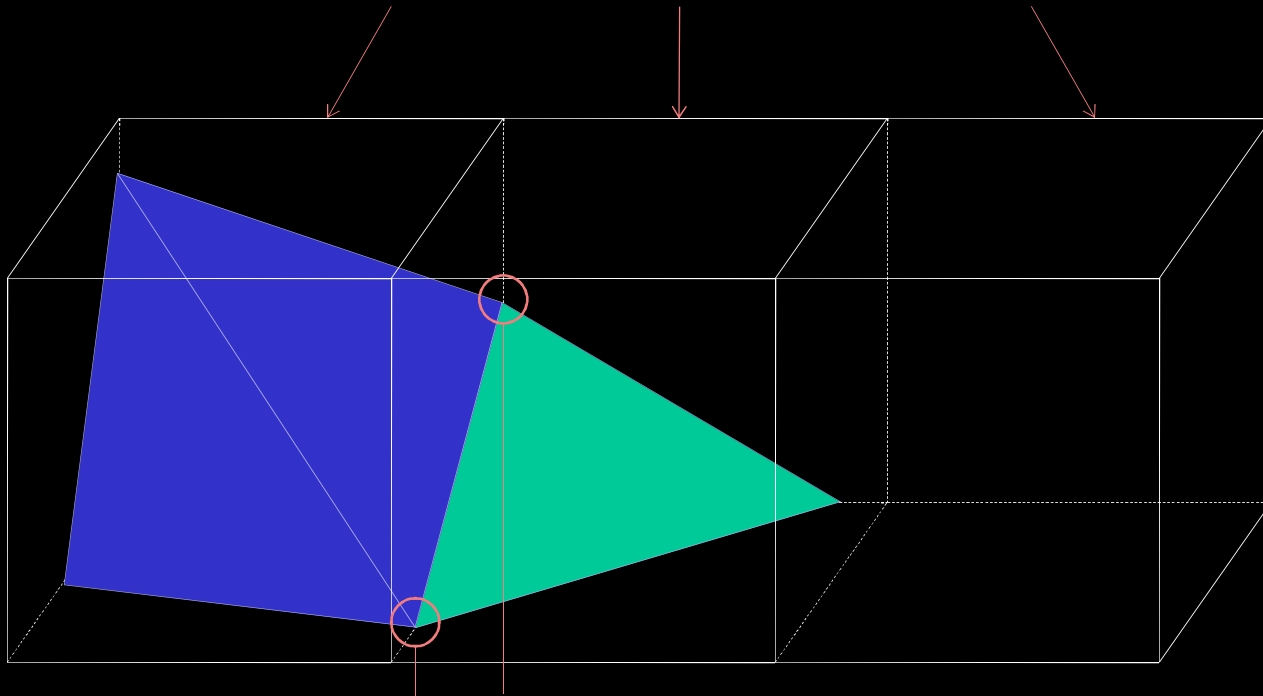
```
DAX_WORKLET void CellGradient(...)
{
    dax::exec::Cell cell(work);
    dax::Vector3 parametric_cell_center
        = dax::make_Vector3(0.5, 0.5, 0.5);

    dax::Vector3 value = cell.Derivative(
        parametric_cell_center,
        points,
        point_attribute,
        0);
    cell_attribute.Set(work, value);
}
```

# Threaded Programming is Hard

## Example: Marching Cubes

Easy because cubes can be processed in parallel, right?



How do you resolve coincident points?  
How do you capture topological connections?

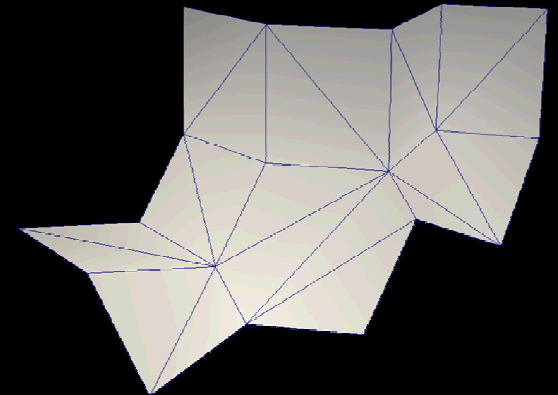
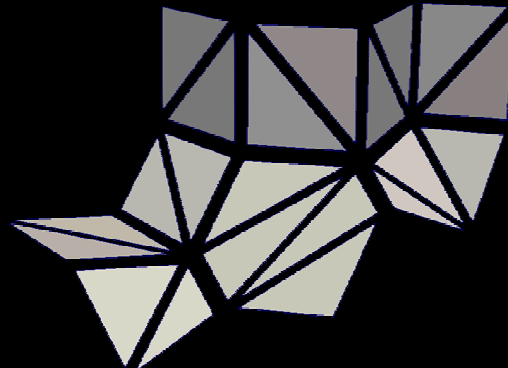
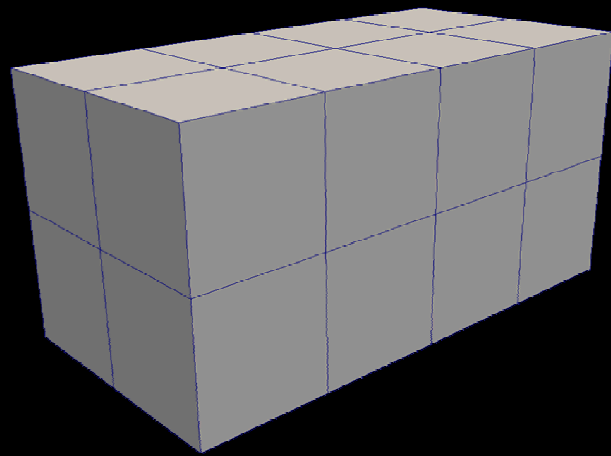
How do you pack the results?

# Work in Progress: Connectivity-Based Algorithms

Input Volume

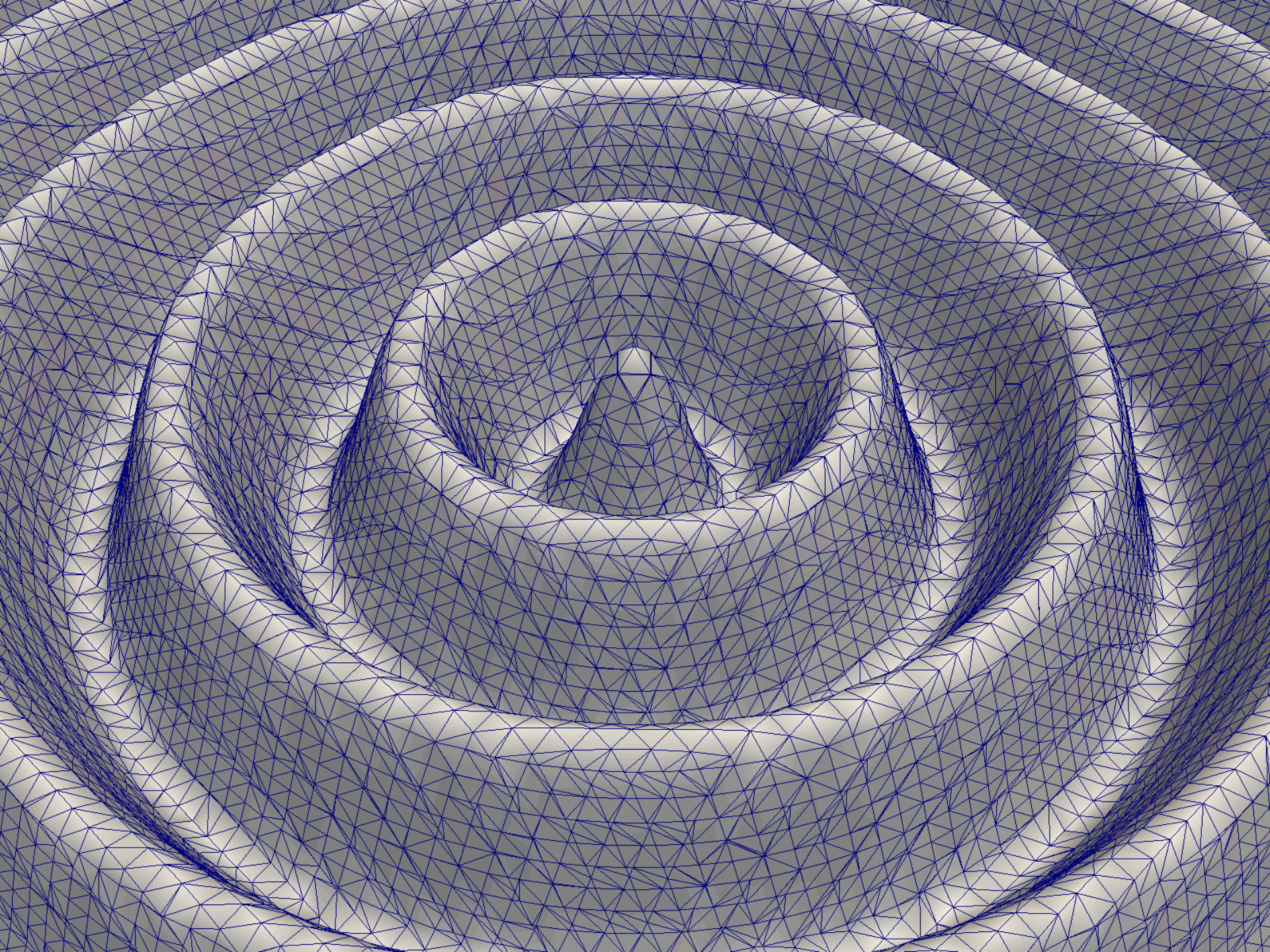
Triangle Soup

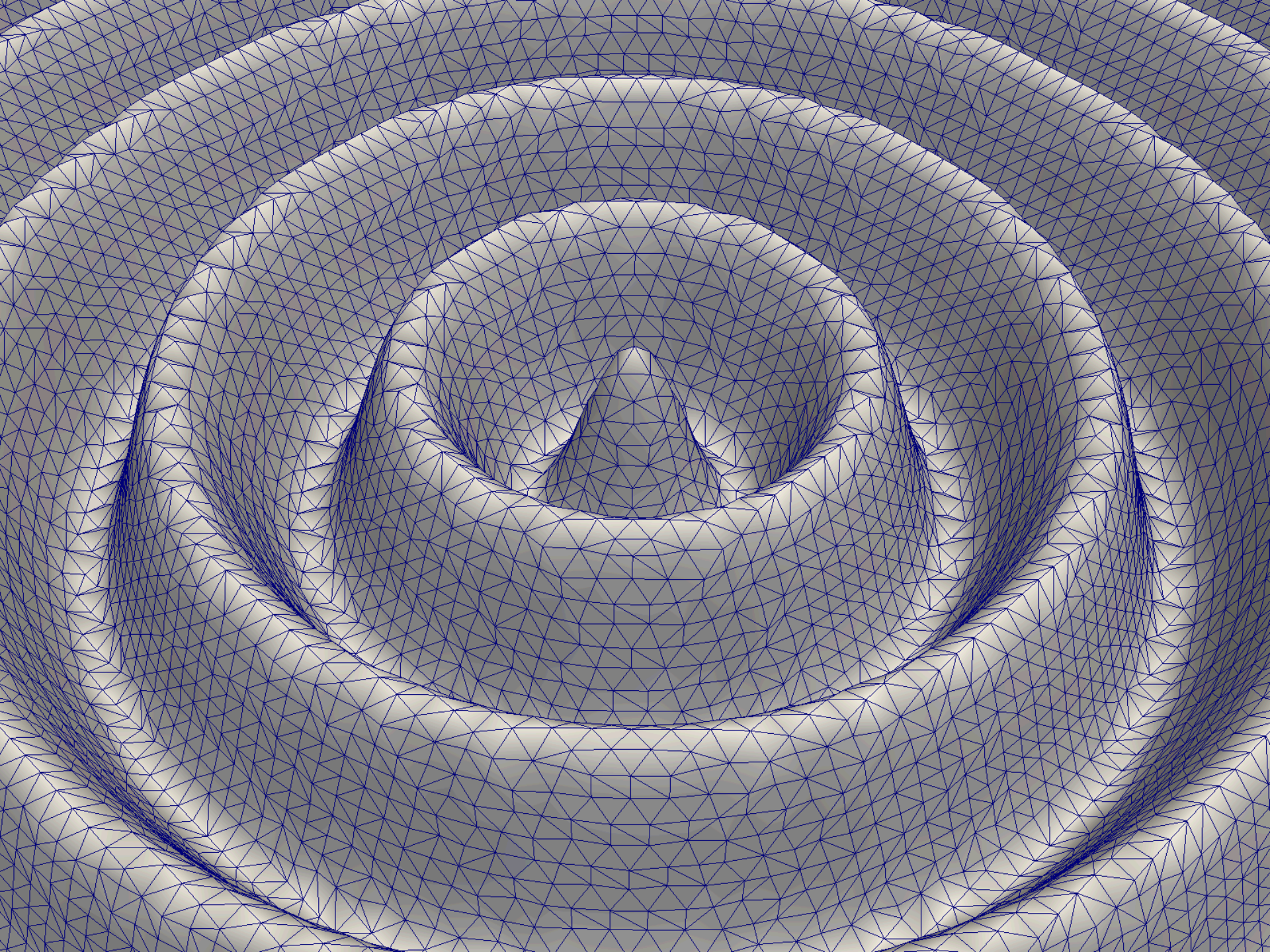
Manifold Surface



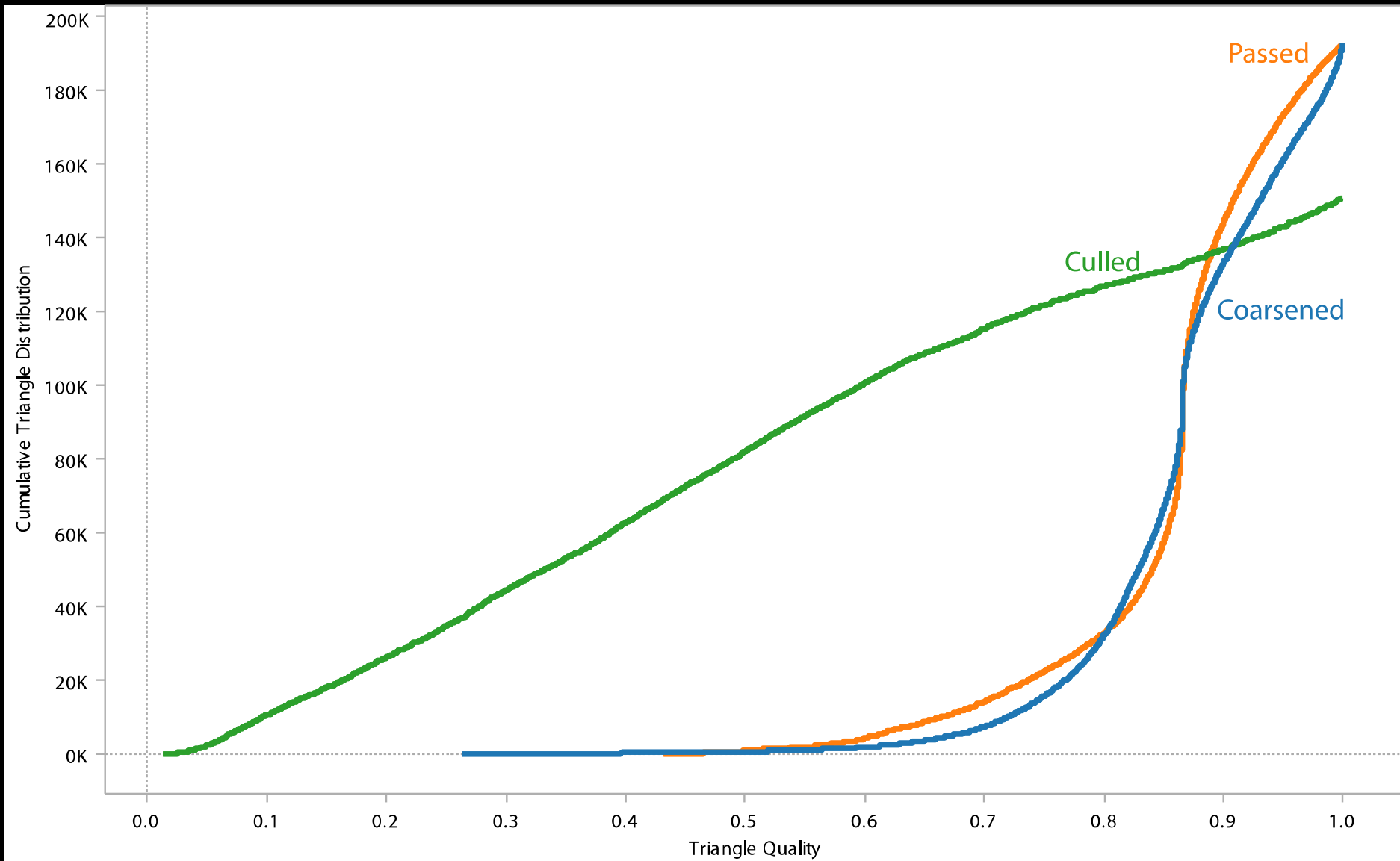
Generate Triangles  
(Local)

Connect Surface  
(Communicative)

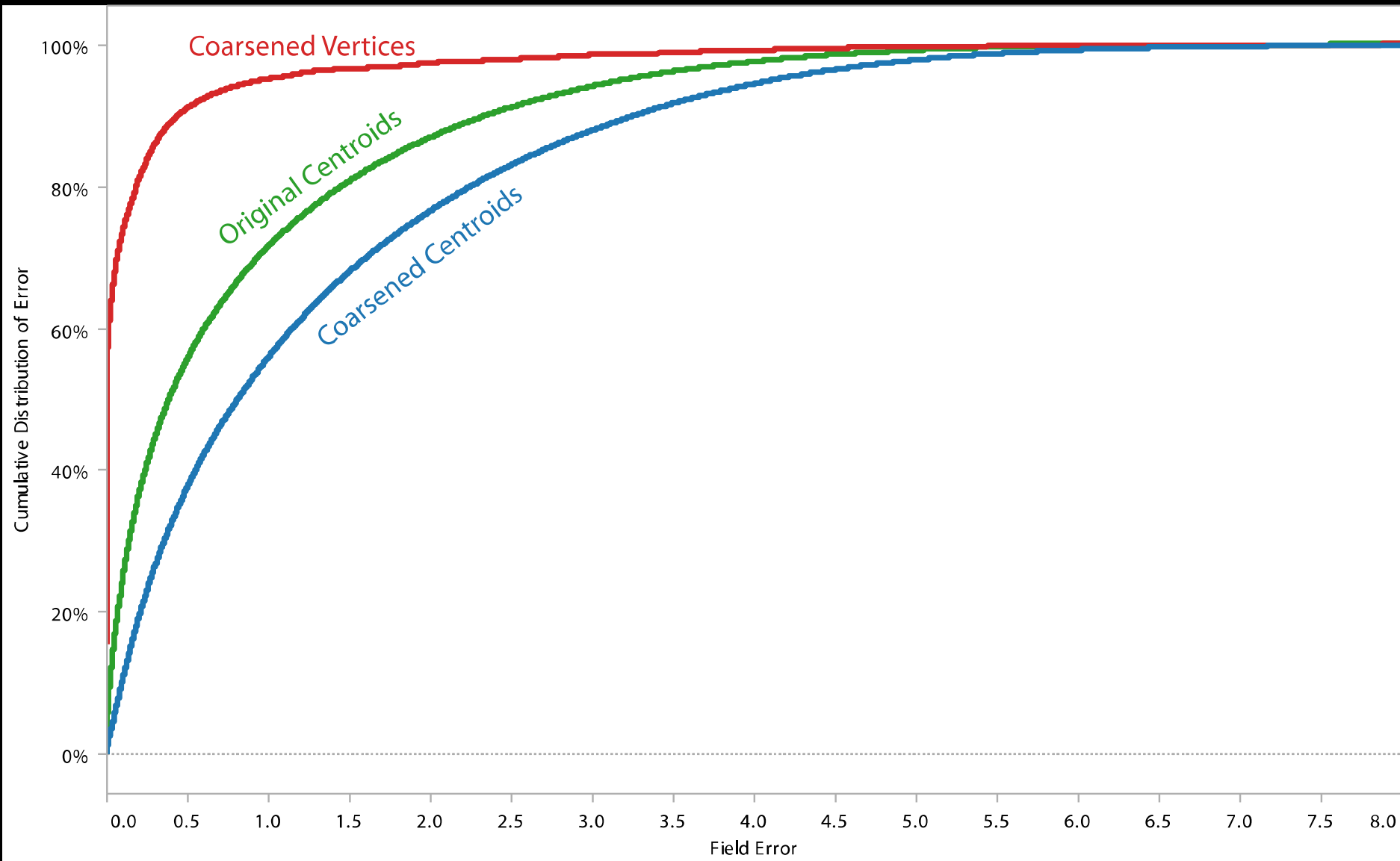




# Coarsened Triangle Quality



# Coarsened Mesh Error



# 1990's

Before consumer market impacted graphics. Single memory/multipipe machines (SGI)

# 2000's

- 2005: Specialized Vis cluster - Distributed memory, commodity clusters tightly coupled with compute platform.
  - Specialized HW (graphics cards, I/O, memory)
- 2010: Distributed memory capability clusters – 'Running on the platform'
  - Highly constrained memory, no specialized HW

# 2020's

- Constrained by power
  - Need: Coupled analysis
- Exascale breaks everything
  - New programming models
  - Billion-way parallelism



## RedRoSE/BlackRoSE



Cielo

## Exascale Platform



Questions?