# Enabling Scalable Genomics Research across Desktop and the Cloud
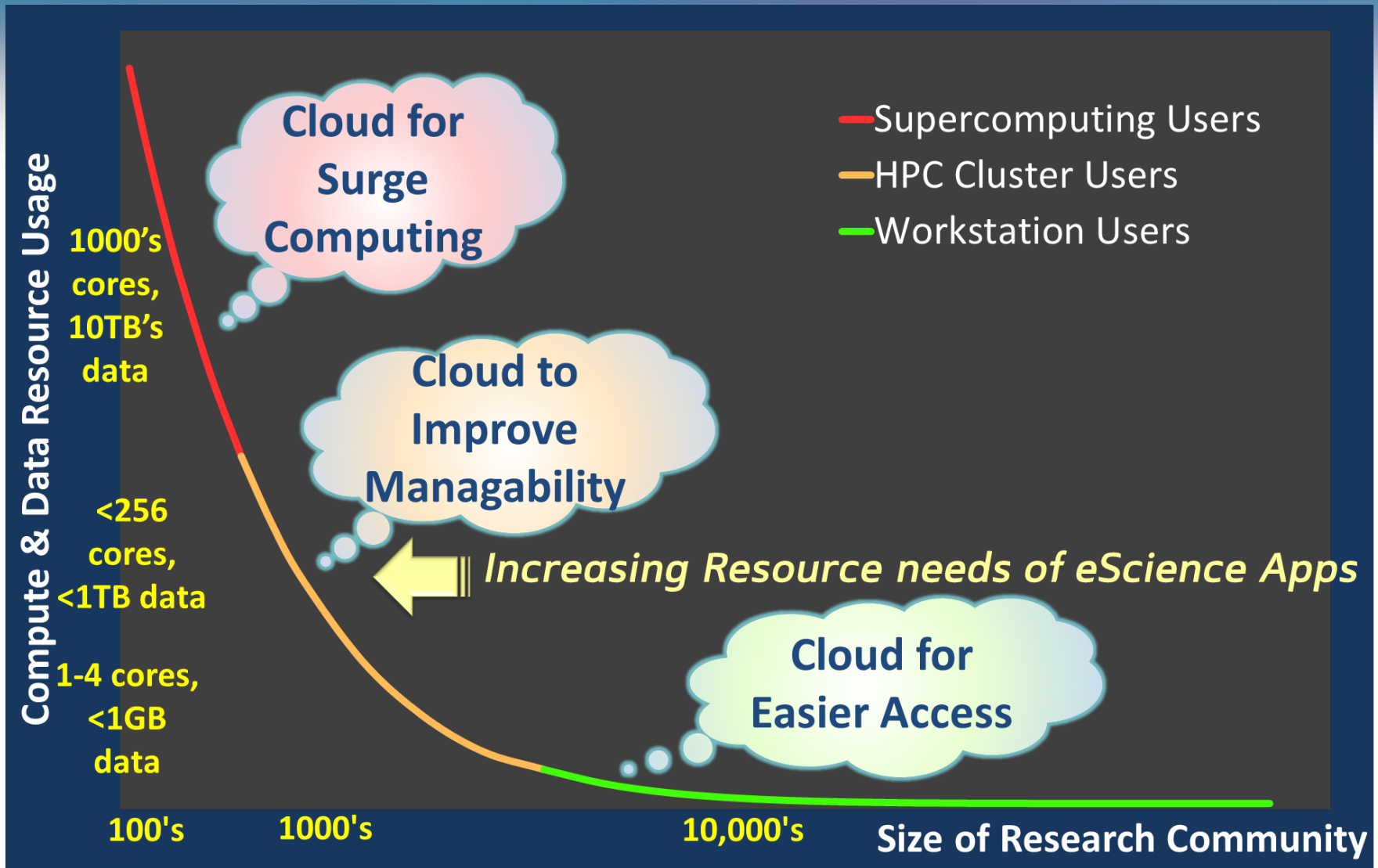
**Yogesh Simmhan**  eScience Group | MSR

Girish Subramanian    DDE Lab | Indiana University
Jennifer Listgarten, David Heckerman, Simon Mercer
& Michael Zyskowski   External Research    | MSR
Wei Lu                 Cloud Comp. Futures | MSR

eScience

# Increasing eScience Resource Needs

- Researchers confronted with ***surfeit of scientific data***
  - Sensors, shared instruments, simulations
- ***Resource needs*** pushed to the next level
  - Collect, process, analyze, visualize data
- ***Variable workloads*** over time
  - Field campaigns, human subject studies
- Distributed ***Collaborations***
  - Share, track, archive, reuse data

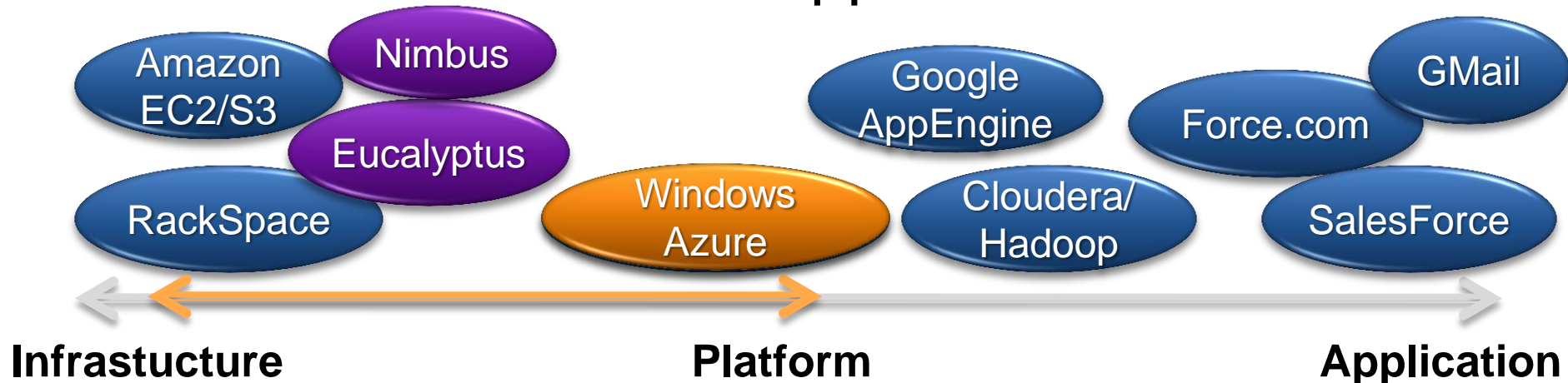# Cloud Opportunities for Science Users

# Cloud Computing Advantages

- On-demand *availability* of computation & storage
- *Scalable* resources from 1 – 100's of nodes
- Ease of resource *management*
- *Economical,* Pay as you go, economies of scale, resources upgraded
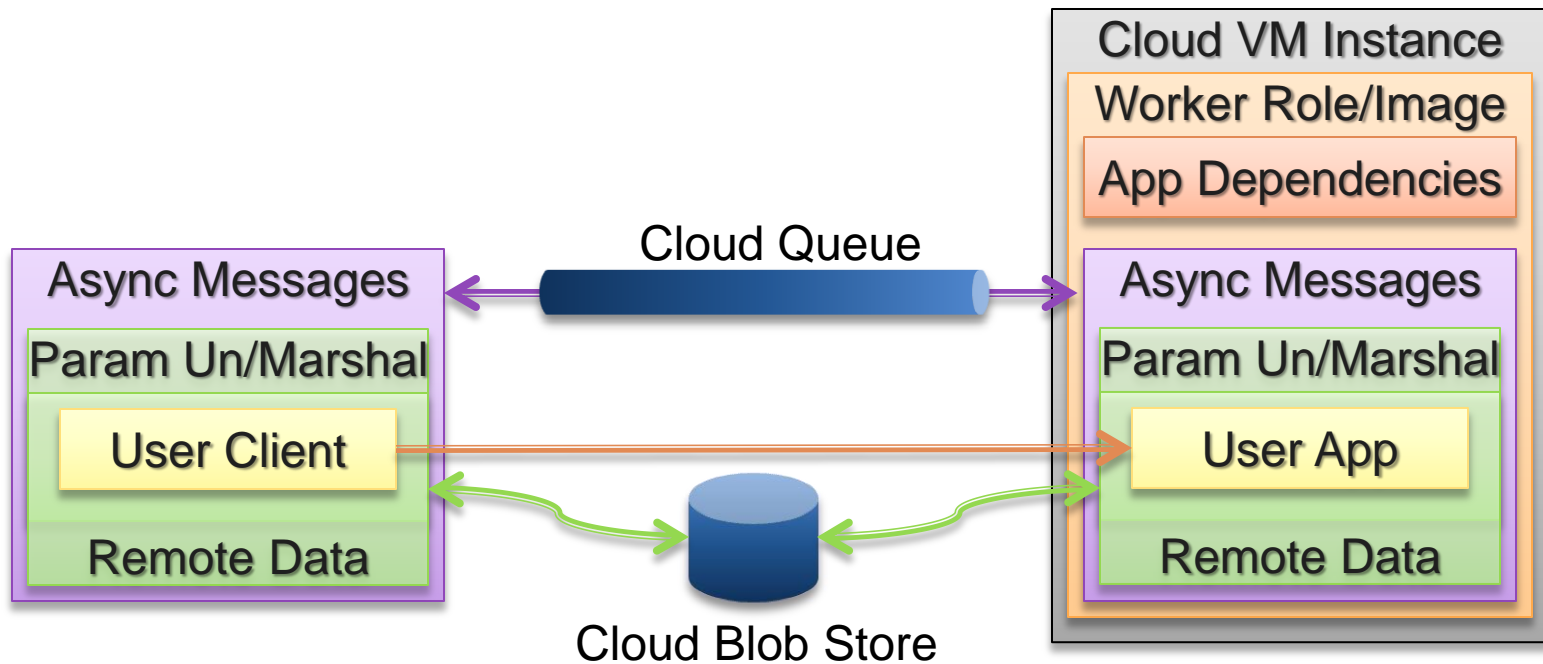- *Simple* interface using REST Web service

# Common Cloud Architectures

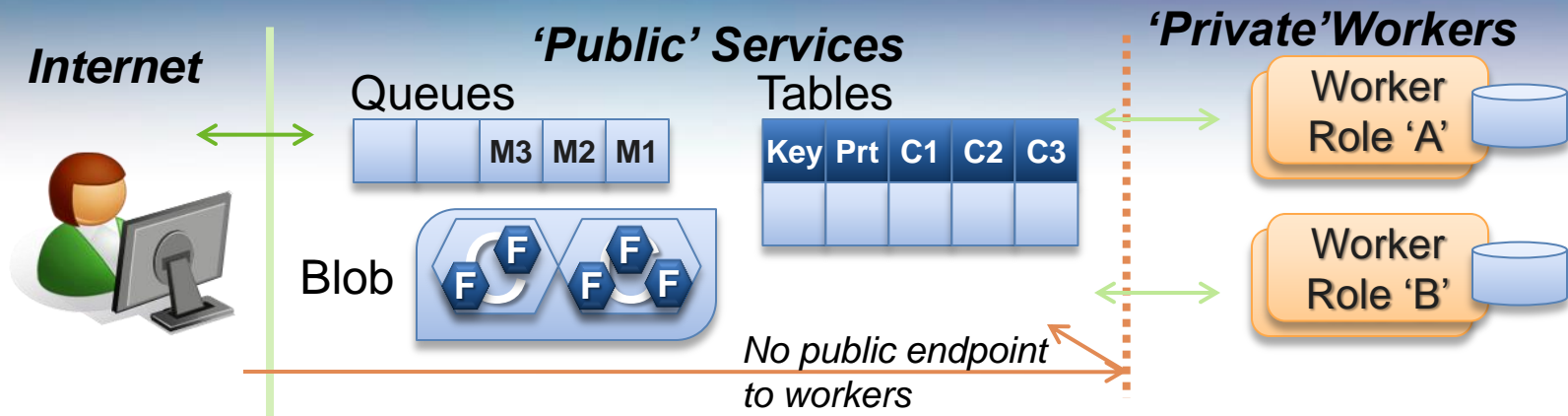- Infrastructure as a Service (*IaaS*)
  - Virtualizes hardware as scalable services
- Platform as a Service (*PaaS*)
  - Provides scalable development platforms
- Software as a Service (*SaaS*)
  - Customize scalable applications

Amazon EC2/S3

Nimbus

Eucalyptus

RackSpace

Windows Azure

Google AppEngine

Cloudera/ Hadoop

GMail

Force.com

SalesForce

**Infrastucture**　　　　　**Platform**　　　　　**Application**

# Challenges for Science Apps across Desktop & Cloud

- Migrating desktop apps to cloud, in part or full
- **Application Deployment**
  - Manage VM images, roles, dependencies
  - Redeploy upon any change
- **Application Execution**
  - Rewrite apps & clients for async remote exec.
  - Param passing, queues, service wrappers
- **Data Access**
  - Access to local files by client & cloud apps

# Migrating Apps to the Azure Cloud

**Internet**

**'Public' Services**

**'Private' Workers**

Queues

| | | M3 | M2 | M1 |
|---|---|---|---|---|

Blob

Tables

| Key | Prt | C1 | C2 | C3 |
|---|---|---|---|---|
| | | | | |

*No public endpoint to workers*

Worker Role 'A'

Worker Role 'B'

## Cloud VM Instance

**Worker Role/Image**

**App Dependencies**

Cloud Queue

| Async Messages | Async Messages |
|---|---|
| Param Un/Marshal | Param Un/Marshal |
| User Client | User App |
| Remote Data | Remote Data |

Cloud Blob Store

# Goal

*Improve accessibility to Cloud:*

- **Deploy** Apps from desktop to cloud
- **Invoke** Cloud Apps from desktop clients and workflows
- Support efficient and automated **file access** across desktop and cloud

*with minimal user overhead and code change*

## » *Generic Worker Framework for Azure*

Common worker role capable of dynamic deployment & execution of registered application from simple desktop commandline clients, workflows or APIs

# Application Deployment
## *Register App with Generic Worker*

- ***Register*** application to make them available for on-demand deployment
- ***Self-contained*** .NET, .exe, Java apps
- Commandline tool: ***register.exe***
- Pass location of ***bin*** directory with app files and unique dependencies
- Pass "***runtime type***" of application for shared dependencies
- Specify ***method*** to run or .exe signature

# Application Deployment
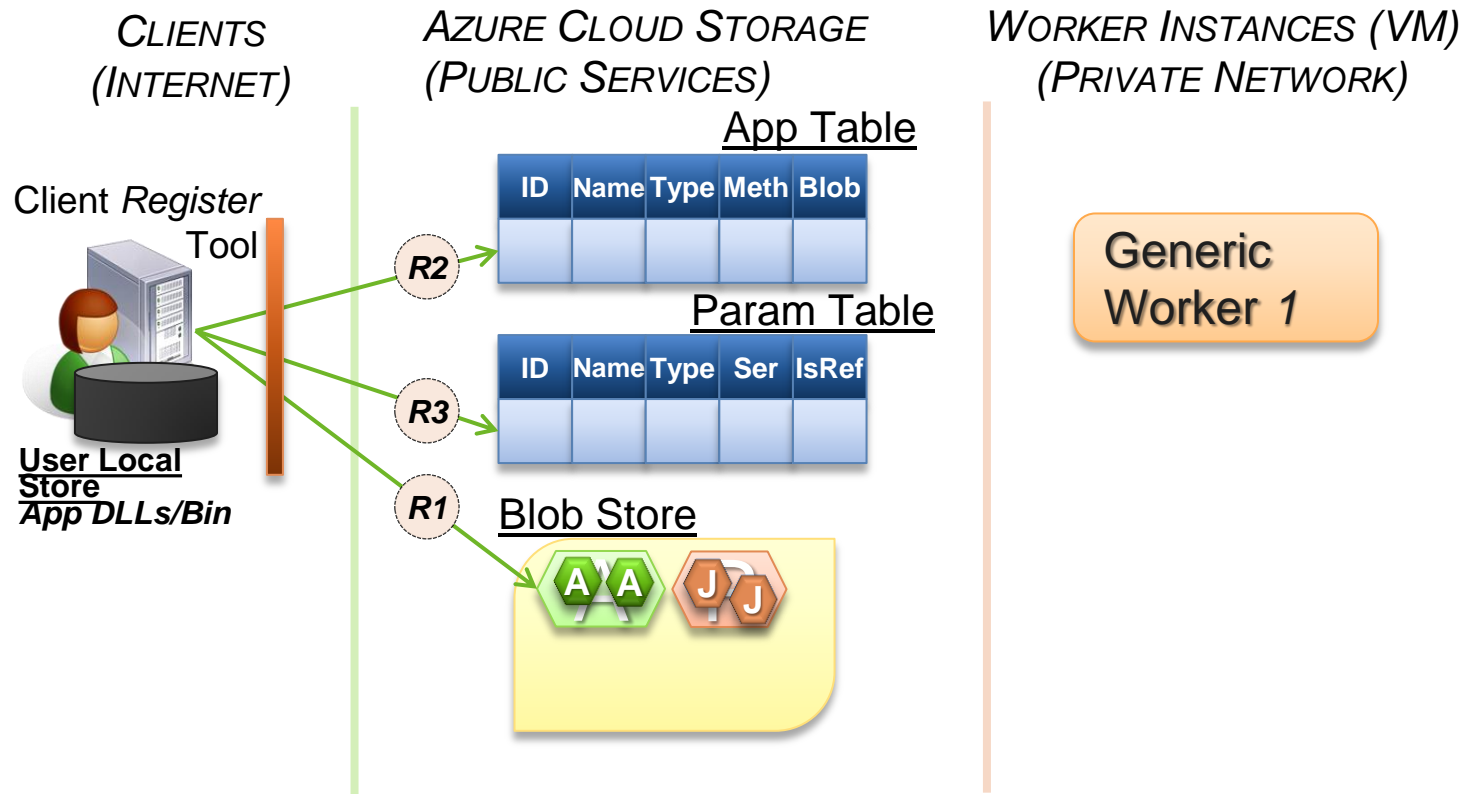## *Sample App Register Command*

```
class MathOps {
    int Add(int i, int j) { return i + j; }
    int Mult(int i, int j) { return i * j; }
}
```

*(a)*

```
> register -type .net35 -name MyMathOps
      -class SampleCloudApp.MathOps
      –appDir c:\SampleCloudApp\bin
> register -type bin -name MyBlastAll
      -cmd "blastall –p blastn –d refseq_rna
-i {1} -o {2}"
      -in #1 file  #2 string  -out #2 file
      –appDir c:\ncbi\blast\bin
```

*(b)*

# Application Deployment
## *App Register Operations*

# Application Execution

- Allow desktop clients to run cloud apps
  - Commandline tool for invocation: *invoke.exe*
  - Pass app input parameters on commanline
  - Internal *Invoker* .NET library to marshal params into XML message, Put on Job Request queue, Poll for response
  - Output XML Message unmarshalled as .NET objects
- Multiple apps can run concurrently in a single worker

# Application Execution
## Sample App Invoke Command
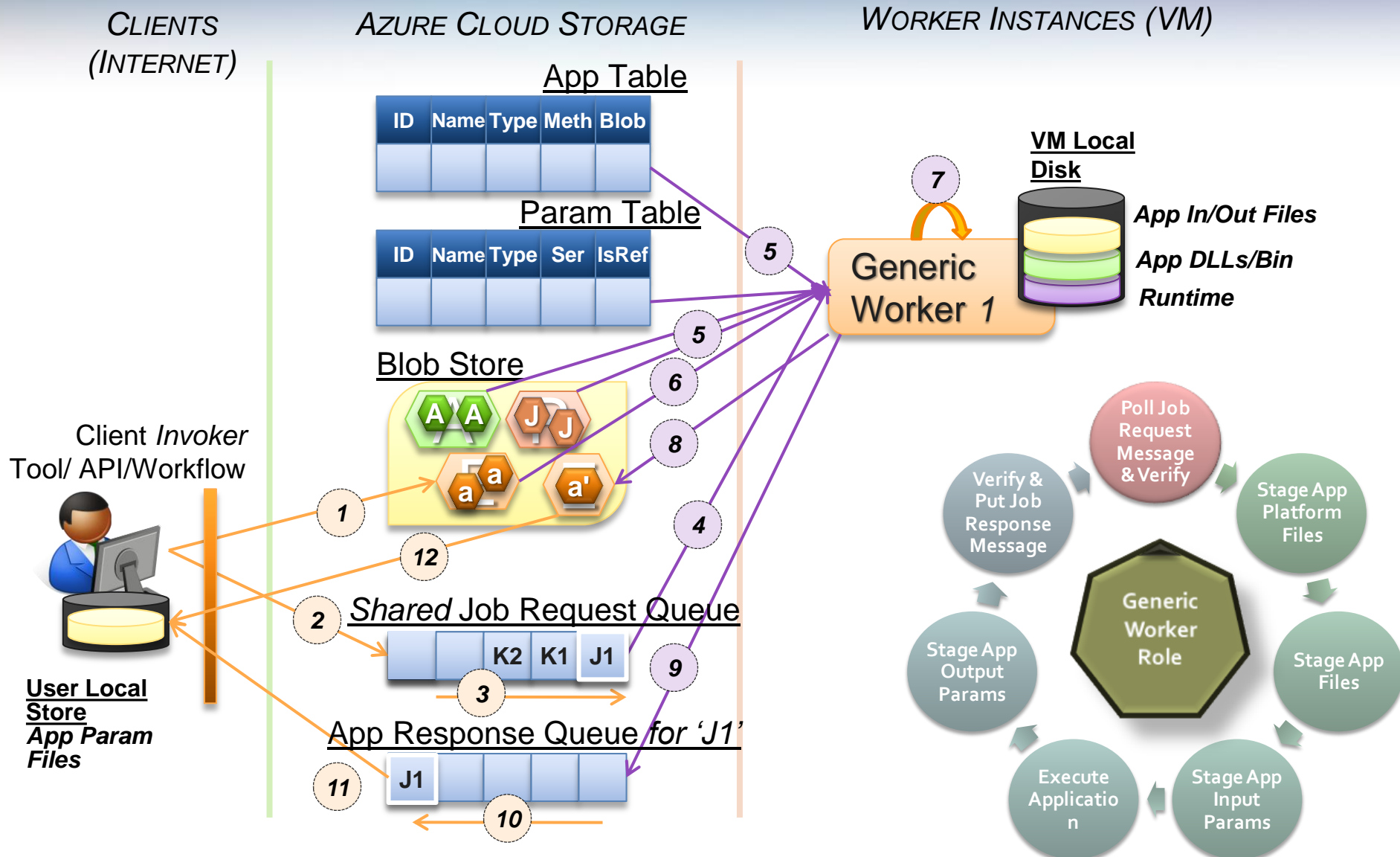
*(a)*

```
> invoke MyMathOps.Add 1 5
Return value: 6
> invoke MyBlastAll input.fasta output.txt
Return value: c:\workdir-036\output.txt
Download Console.Out file (y/n)?
```

*(b)*

```
// int s = (new MathOps()).Add(1,5);
Invoker invkr = new Invoker("MyMathOps);
int s = (int)invkr.Invoke("Add",new[]{1,5});
```

# Application Execution
# *App Invoke Operations*



**CLIENTS (INTERNET)**

**AZURE CLOUD STORAGE**

**WORKER INSTANCES (VM)**

App Table

| ID | Name | Type | Meth | Blob |
|----|------|------|------|------|
|    |      |      |      |      |

Param Table

| ID | Name | Type | Ser | IsRef |
|----|------|------|-----|-------|
|    |      |      |     |       |

Blob Store

Client *Invoker* Tool/ API/Workflow

**User Local Store**
*App Param Files*

*Shared* Job Request Queue

App Response Queue *for 'J1'*

VM Local Disk

Generic Worker 1

*App In/Out Files*

*App DLLs/Bin*

*Runtime*

Poll Job Request Message & Verify

Stage App Platform Files

Verify & Put Job Response Message

Generic Worker Role

Stage App Files

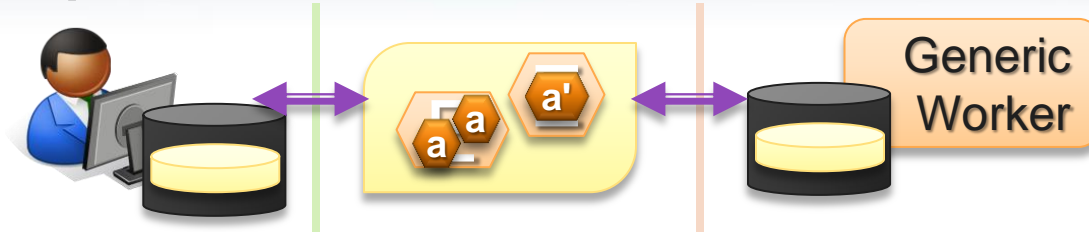Stage App Output Params

Stage App Input Params

Execute Application

# Data Access
## *Automatic & Efficient Data Transfers*

- Desktop:Worker file transfer thru Blob Store



- Generic Worker recognizes "file" params during registration

- Framework does automatic JIT data transfer
  - Clients & apps continue to pass local file paths

- File "references" ensure transfer if needed

- Basic caching to reduce cloud:local transfers

```
> invoke MyBlastAll input.fasta output.txt
Return value: c:\workdir-036\output.txt
Download Console.Out file (y/n)?
```
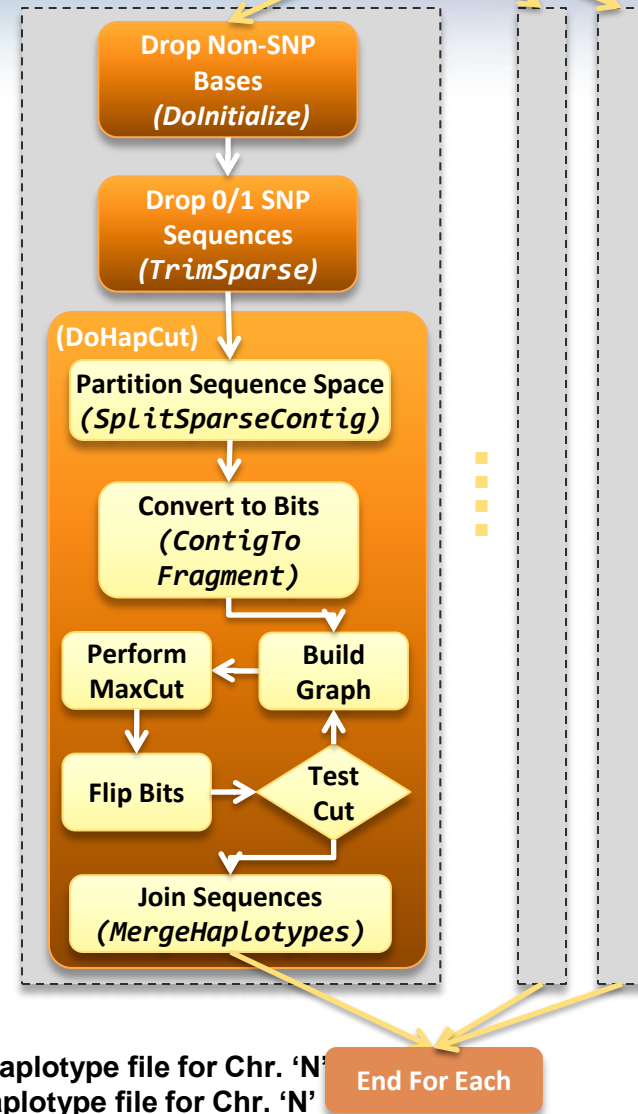*(a)*

# Genome Phasing

- Microsoft Biology Foundation
  - Library of genomics algorithms, data structures & file parsers
- Genome phasing separates parent haplotype sequences from sequence reads
- Phasing workflow includes compute intensive algorithm & pre-, post- procesors
- Achieve parallelism across chromosomes
- Goal: Orchestrate workflow and pre-, post- activities locally but ship compute heavy to cloud

*With David Heckerman, Simon Mercer & Michael Zyskowski*

➢ Aligned Sequence File for Chr. 'N'
➢ Reference SNP DB File
➢ Chromosome Number 'N'

For Each Chromosome 'N' Do in Parallel

Drop Non-SNP Bases
*(DoInitialize)*

Drop 0/1 SNP Sequences
*(TrimSparse)*

(DoHapCut)

Partition Sequence Space
*(SplitSparseContig)*

Convert to Bits
*(ContigTo Fragment)*

Perform MaxCut | Build Graph

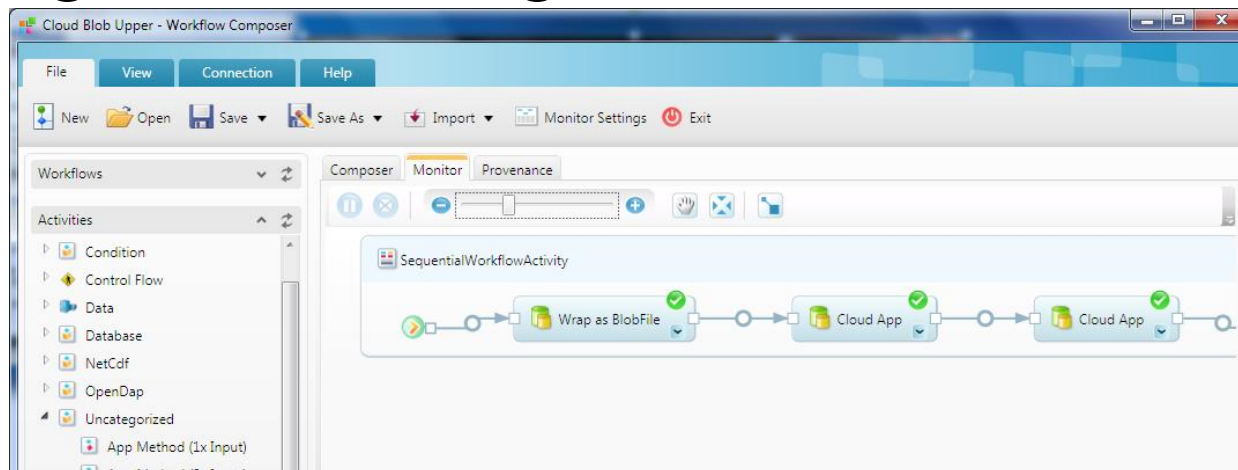Flip Bits | Test Cut

Join Sequences
*(MergeHaplotypes)*

➢ 'Mom' Haplotype file for Chr. 'N'
➢ 'Dad' Haplotype file for Chr. 'N'

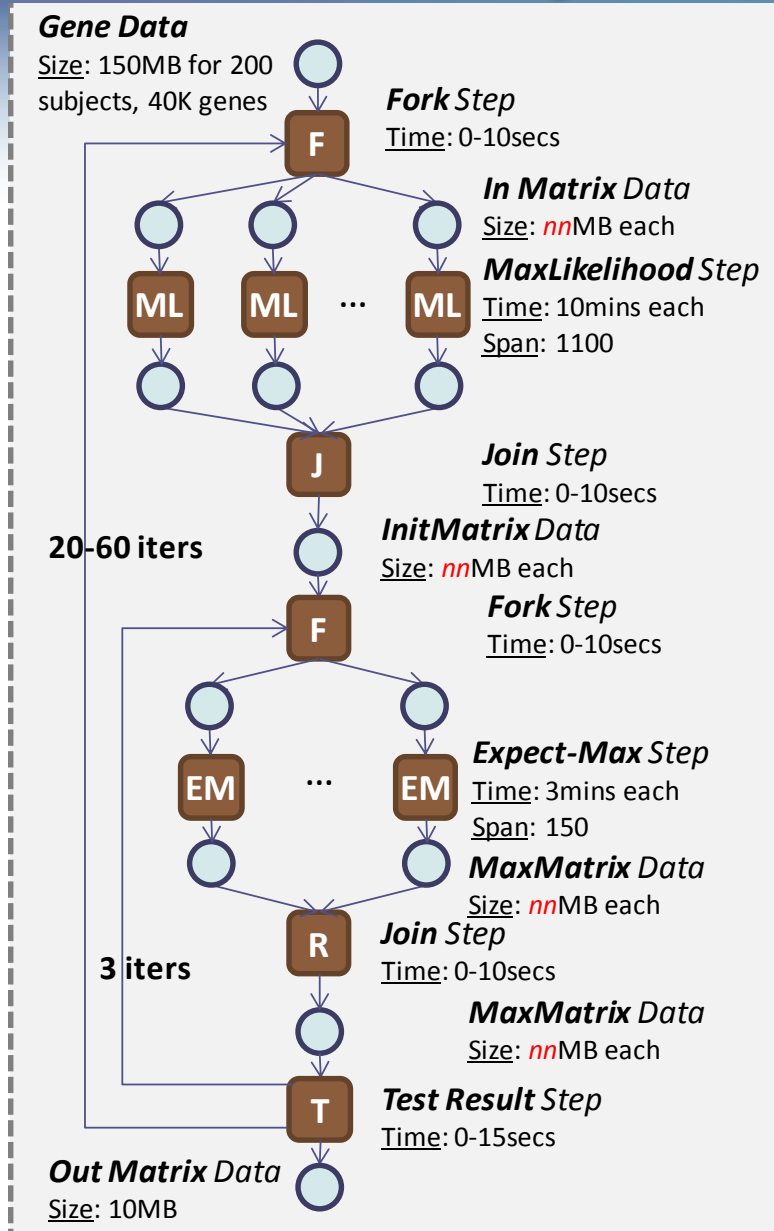End For Each

# Genome Phasing Workflow

- Trident as workflow environment
- *CloudApp* activity can invoke any registered app in Generic Worker using API
- Data transfers are transparent
- Potential to host a suite of composable MBF genomics algorithms in Cloud
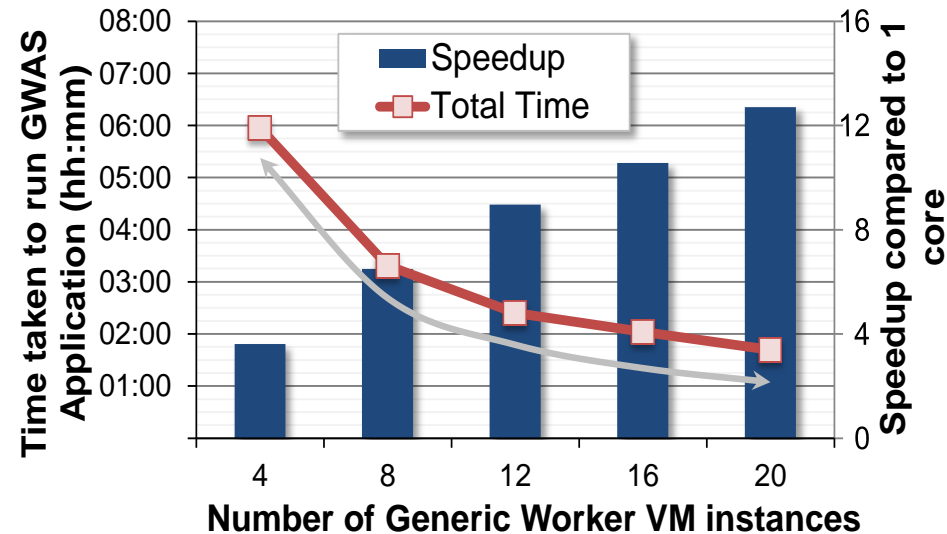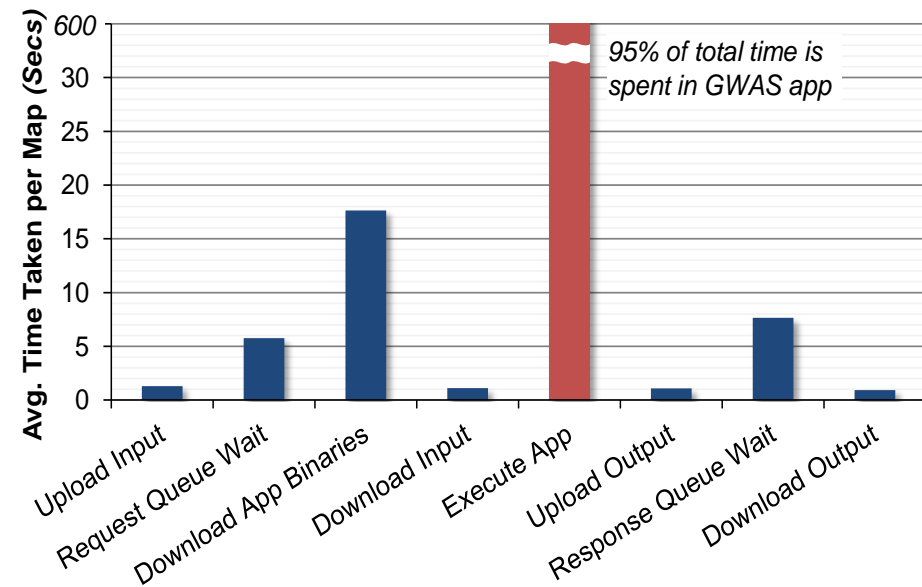
# Genome Wide Association Study

- Iterative MapReduce of ML & EM methods

- Compute intensive
  - ~18hrs on 300 HPC cores

- .NET code orchestrates execution, batch submission, file & parameter passing

- Need to share with research community

*With Jennifer Listgarten & David Heckerman*



**Gene Data**
Size: 150MB for 200 subjects, 40K genes

**Fork** Step
Time: 0-10secs

**In Matrix** Data
Size: *nn*MB each

**MaxLikelihood** Step
Time: 10mins each
Span: 1100

**Join** Step
Time: 0-10secs

**InitMatrix** Data
Size: *nn*MB each

**Fork** Step
Time: 0-10secs

20-60 iters

**Expect-Max** Step
Time: 3mins each
Span: 150

**MaxMatrix** Data
Size: *nn*MB each

**Join** Step
Time: 0-10secs

3 iters

**MaxMatrix** Data
Size: *nn*MB each

**Test Result** Step
Time: 0-15secs

**Out Matrix** Data
Size: 10MB

# GWAS application on Azure

- ML and EM methods registered as apps
- Fork-Join to Generic Workers replaced HPC parametric sweep batch submission



*95% of total time is spent in GWAS app*

- Allow easy switching between HPC and Cloud environments

# Future Work

- Support further application runtime types
  - Beyond .NET & Exe to Java, MatLab, etc.
- Dynamic worker instance scale-out based on collective demand
- Leverage data locality in worker file cache
  - E.g. workflow data pipelines
- Metrics for automated cross-scheduling workflows across Desktop, HPC and Cloud

# Conclusions

- Cloud has opportunities to democratize science, but has to be accessible

- Reducing overhead to migrating existing desktop apps to Cloud can drive adoption

- App deployment, remote execution and file handling often require code change

- Generic Worker reduces user and coding overhead with limited performance penalty

- Intuitive execution across desktop & Cloud from commandline, workflows & APIs