

Hands-on 3: Local and Distributed Execution

Daniel Crawl

Shweta Purawat

Goals

- Create an actor to run BLAST in different environments
 - Local execution (blastall)
 - (Part 2) Distributed execution (mpiblast) on a cluster
- Configure ExecutionChoice to run BLAST
 - Inputs, Outputs, Parameters
 - Local and MPI choices

Running BLAST Locally

- Example command line:

```
blastall -i query -d ref -p blastn -m 8 -e 1E-5 > align
```

Running BLAST Locally

- Example command line:

```
blastall -i query -d ref -p blastn -m 8 -e 1E-5 > align
```

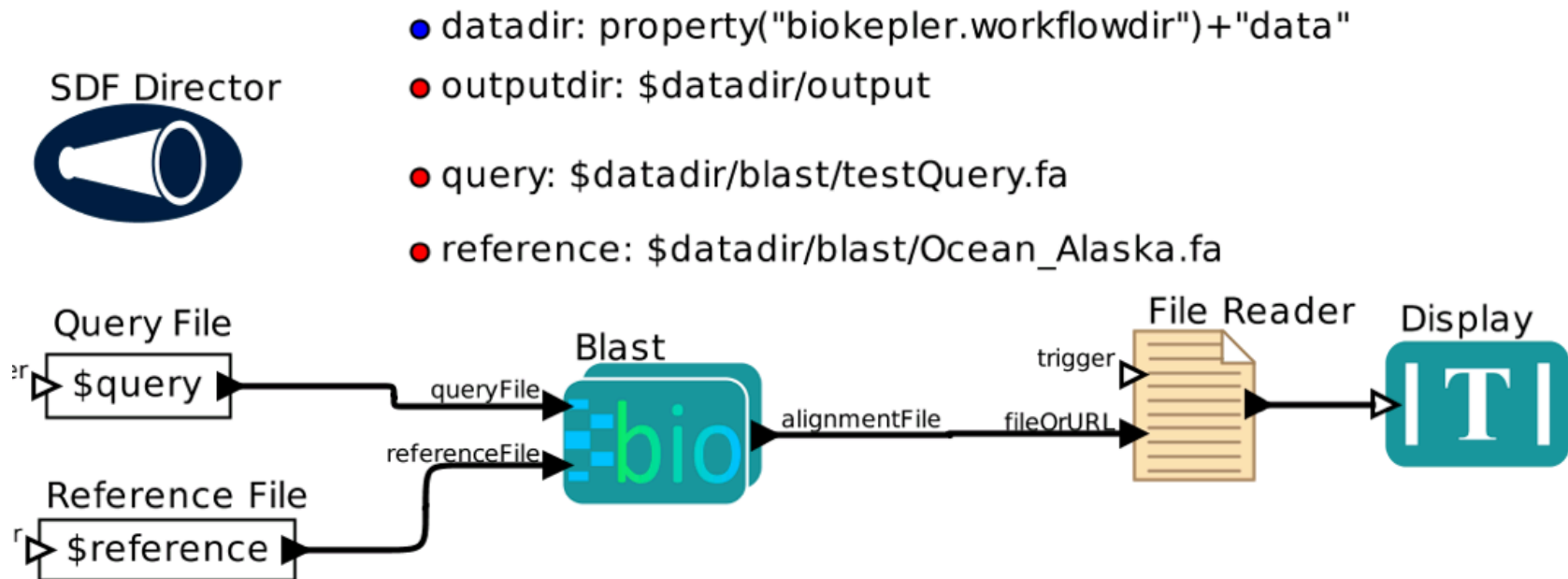
Inputs

Parameters

Output

Step 1: Open Blast workflow

- Workflow is /home/biokepler/Blast.kar



Step 2: Set Parameter Values

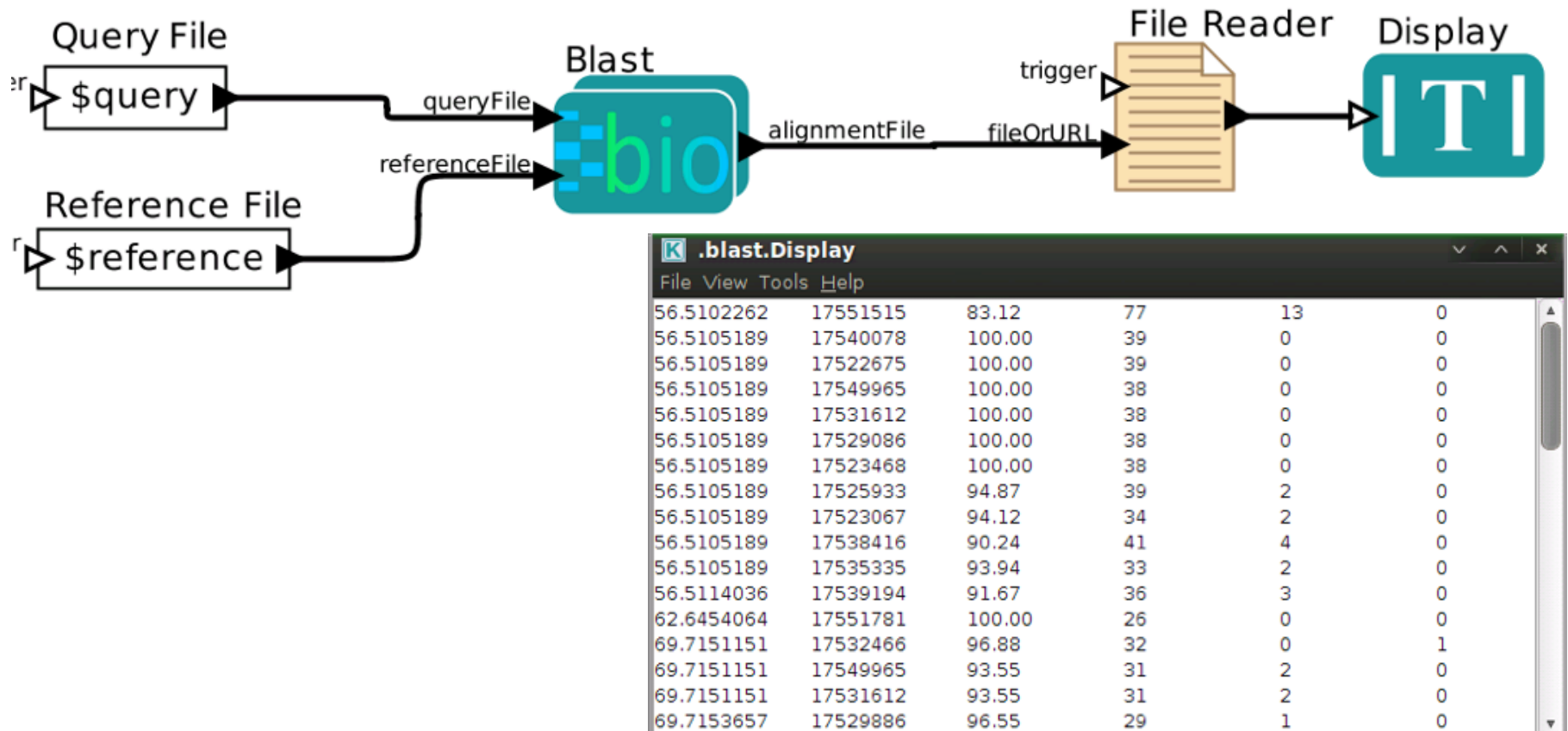
| | |
|------------------------|---|
| program: | <input type="text" value="blastall"/> |
| Input File Parameters | |
| queryFile (-i): | <input type="text"/> |
| referenceFile (-d): | <input type="text"/> |
| Output File Parameters | |
| checkOutputTimestamp: | <input checked="" type="checkbox"/> |
| alignmentFile (>): | <input type="text" value="\$HOME/alignment.txt"/> |
| Parameters | |
| additionalOptions: | <input type="text" value="-e 1E-5 -p blastn -m 8"/> |

Step 3: Run Finished Workflow

SDF Director



- datadir: `property("biokepler.workflowdir")+"data"`
- outputdir: `$datadir/output`
- query: `$datadir/blast/testQuery.fa`
- reference: `$datadir/blast/Ocean_Alaska.fa`



mpiBLAST

- mpiBLAST is open source MPI based implementation of database segmentation for parallel BLAST search
- Super-linear performance gain with database segmenting technique
- Ideal database fragment is largest fragment that can sit in the memory
- Making fragments smaller than available memory adds to the overhead

mpiBLAST Algorithm

mpiBLAST algorithm consists of three steps:

- Segmenting and Distributing the database
- Running mpiBLAST queries on each node
- Merging the results from each node in a single output file

mpiBlast commands

- **Formatting and Segmenting the database**

`mpiformatdb -n NP -i ref.fa`

- **Querying the database and merging the results**

`mpirun -n NP mpiblast -p blastp -d ref.fa -i query.fas -o blast_results.txt`

Goal: To run mpiBlast on HPC Cluster remotely
through Kepler workflow: SDSC Gordon
Supercomputer

```
mpirun -n NP mpiblast -p blastp -i query -d ref
```

```
-m 8 -e 1E-5 > align
```

Inputs

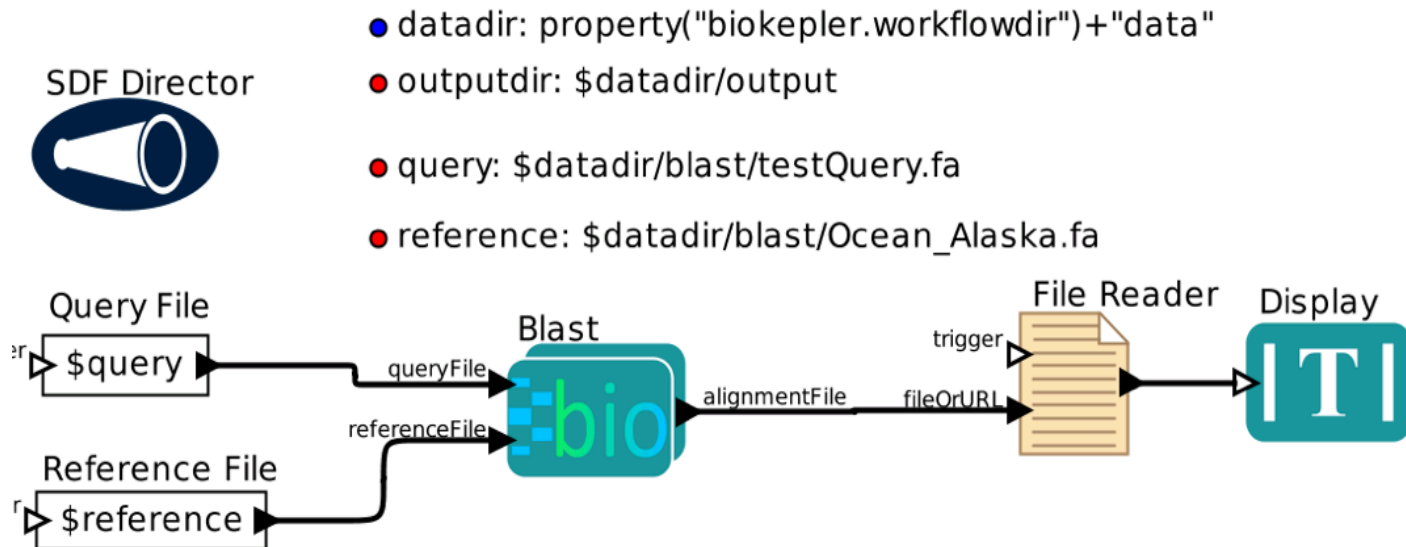
Parameters

Output

NP – multiple parallel processes

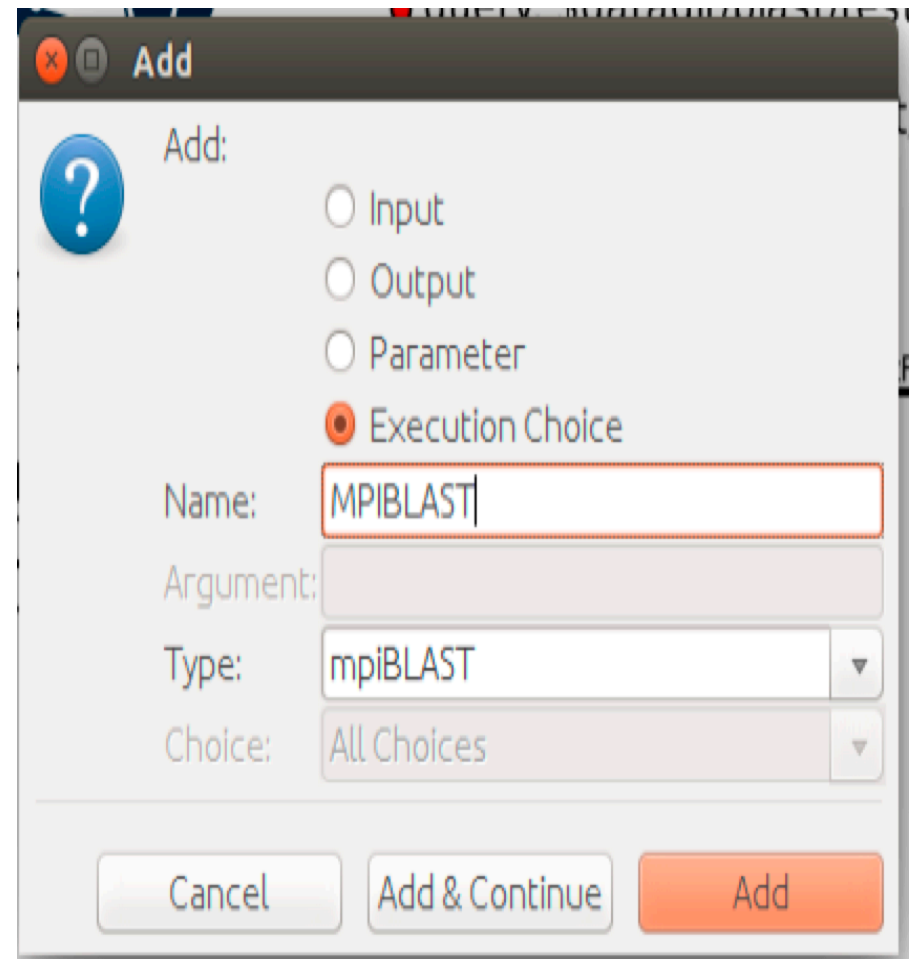
Step 1: Open blast workflow

- Go to Step 2, if Blast.kar is already open
- Else the workflow: /home/biokepler/Blast.kar



Step 2: Add new execution choice for MPIBLAST

- Double-click on the Blast actor and click on Add.
- Change radio button to Execution Choice.
- Select Type: mpiBLAST.
- Set Name to MPIBLAST.
- Click Add to close the dialog



Step 3: Configure the New Execution Choice

- Double-click on Blast actor
- Change to MPIBLAST tab
- Configure MPIBLAST tab:

Please use Gordon username assigned to you in place of **LOGIN**.

First parameter: TargetHost

- Edit “TargetHost” parameter and set its value to be LOGIN@gordon.sdsc.edu

Second parameter: cmdFile

- If already set move next parameter
- Else set its value to be \$HOME/mpiBLAST/blastp_gordon.sh

Third parameter: commandLine

- Keep as it is. Move to next parameter.

Fourth parameter: outputFile

- If already set move next parameter
- Else set its value to be \$alignmentFile

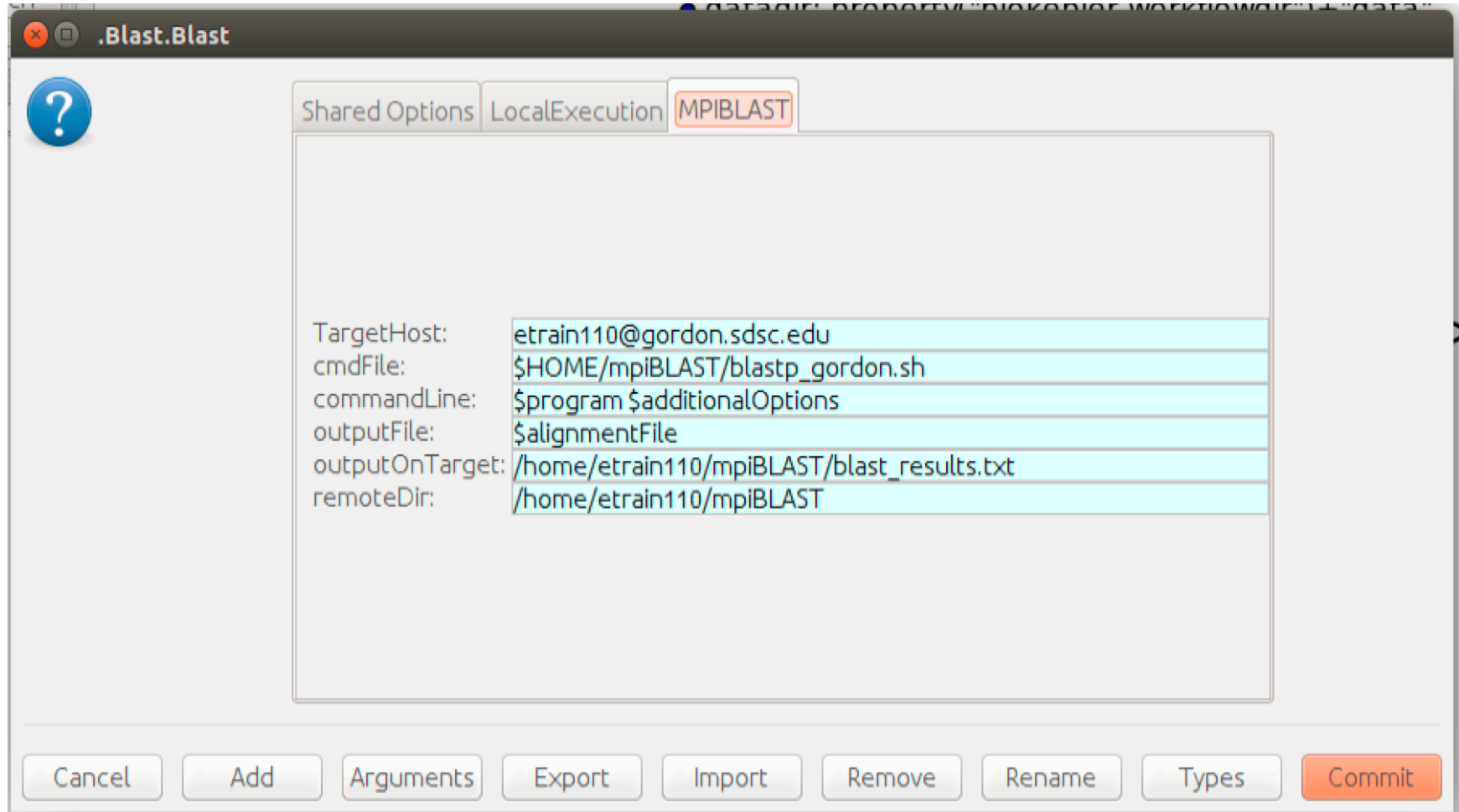
Fifth parameter: outputOnTarget

- Set its value to be /home/LOGIN/mpiBLAST/blast_results.txt

Sixth parameter: remoteDir

- Set its value to be /home/LOGIN/mpiBLAST

Step 3: Configure the New Execution Choice

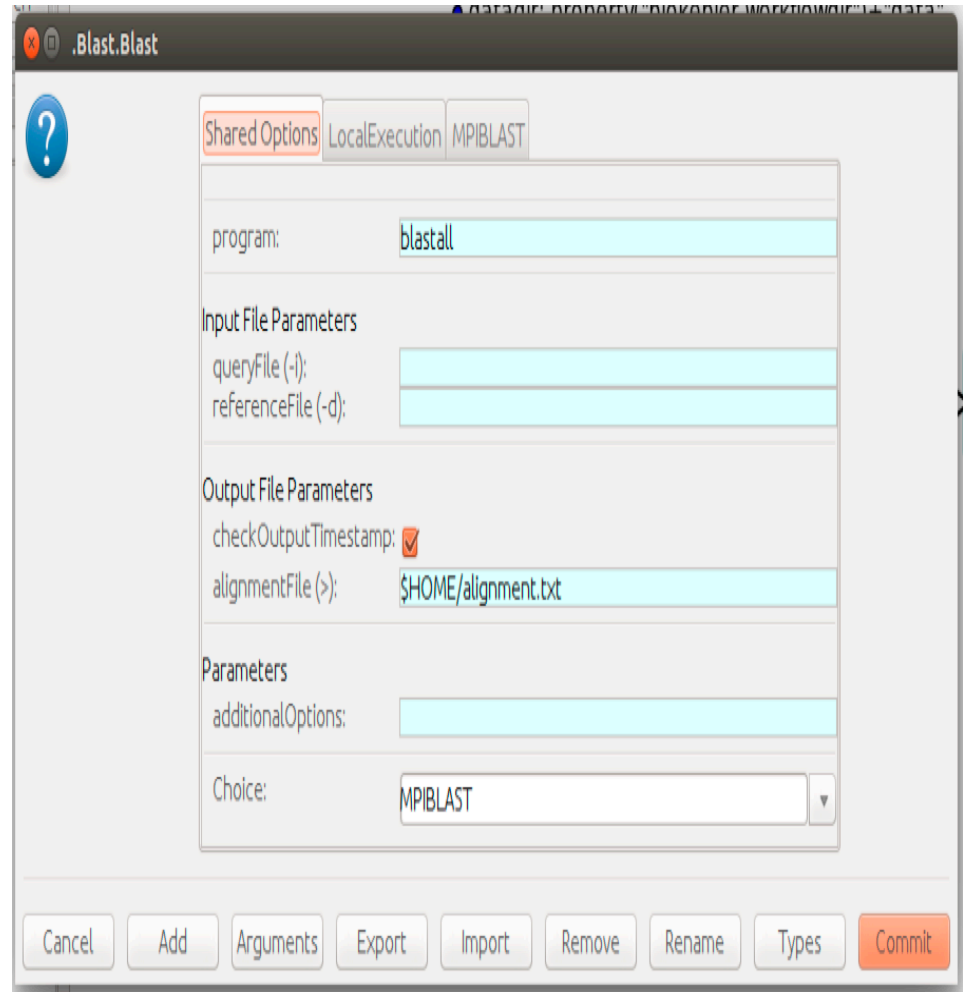


The screenshot shows a window titled ".Blast.Blast" with a help icon (question mark in a blue circle) in the top-left corner. The window has three tabs: "Shared Options", "LocalExecution", and "MPIBLAST". The "MPIBLAST" tab is selected and highlighted with a red border. Below the tabs is a large text area containing configuration fields. At the bottom of the window is a row of buttons: "Cancel", "Add", "Arguments", "Export", "Import", "Remove", "Rename", "Types", and "Commit".

| | |
|-----------------|--|
| TargetHost: | etrain110@gordon.sdsc.edu |
| cmdFile: | \$HOME/mpiBLAST/blastp_gordon.sh |
| commandLine: | \$program \$additionalOptions |
| outputFile: | \$alignmentFile |
| outputOnTarget: | /home/etrain110/mpiBLAST/blast_results.txt |
| remoteDir: | /home/etrain110/mpiBLAST |

Step 4: Set the Execution Choice to MPIBLAST

- Double-click on the Blast actor and click on Add.
- Go to “Shared Options” tab
- Set Choice to MPIBLAST.
- Click on commit to save the changes and close the configuration box.



Step 5: Run Finished Workflow

- Save and Run the workflow.

Workflow

SDF Director

- datadir: property("biokepler.workflowdir")+"data"
- outputdir: \$datadir/output
- query: \$datadir/blast/testQuery.fa
- reference: \$datadir/blast/Ocean_Alaska.fa

File Reader Display

.Blast.Display

File Tools Help

| | | | | | | |
|-----------------------|----------------------------|-------|-----|----|---|-----|
| sp Q4V8X4 ACBD6_DANRE | gi 136628751 gb EBP38214.1 | 34.52 | 84 | 55 | 0 | 192 |
| sp Q4V8X4 ACBD6_DANRE | gi 142867380 gb EDB14451.1 | 35.29 | 85 | 55 | 0 | 179 |
| sp Q4V8X4 ACBD6_DANRE | gi 139289341 gb ECE46419.1 | 37.50 | 88 | 55 | 0 | 192 |
| sp Q4V8X4 ACBD6_DANRE | gi 139710199 gb ECG93695.1 | 41.33 | 75 | 44 | 0 | 201 |
| sp Q4V8X4 ACBD6_DANRE | gi 139710199 gb ECG93695.1 | 41.89 | 74 | 43 | 0 | 203 |
| sp Q4V8X4 ACBD6_DANRE | gi 138897158 gb ECC33175.1 | 37.62 | 101 | 59 | 3 | 202 |
| sp Q4V8X4 ACBD6_DANRE | gi 140000798 gb ECI93000.1 | 35.00 | 100 | 64 | 1 | 180 |
| sp Q4V8X4 ACBD6_DANRE | gi 140139230 gb ECJ80050.1 | 35.56 | 90 | 57 | 1 | 201 |
| sp Q4V8X4 ACBD6_DANRE | gi 142345366 gb ECX41763.1 | 35.29 | 85 | 55 | 0 | 179 |
| sp Q4V8X4 ACBD6_DANRE | gi 140378174 gb ECL36374.1 | 40.54 | 74 | 44 | 0 | 203 |
| sp Q4V8X4 ACBD6_DANRE | gi 141406115 gb ECR60568.1 | 34.02 | 97 | 63 | 1 | 180 |
| sp Q4V8X4 ACBD6_DANRE | gi 138607027 gb ECA76768.1 | 41.89 | 74 | 43 | 0 | 203 |
| sp Q4V8X4 ACBD6_DANRE | gi 138607027 gb ECA76768.1 | 40.54 | 74 | 44 | 0 | 202 |

mpiBLAST Performance

| #Nodes | Run Time (sec) | Speed-Up |
|--------|----------------|----------|
| 1 | 80775 | 1.00 |
| 4 | 8752 | 9.23 |
| 8 | 4548 | 17.76 |
| 16 | 2437 | 33.15 |
| 32 | 1350 | 59.83 |
| 64 | 851 | 94.92 |
| 128 | 474 | 170.41 |

Parallel Computing: Software Technology, Algorithms, Architectures & Applications

Gerhard Joubert, Wolfgang Nagel, Frans Peters, Wolfgang Walter
Elsevier, Sep 23, 2004