

NEAT

Adding modules

Patrick Schorderet

Patrick.schorderet@molbio.mgh.harvard.edu

Jan 2015

1	INTRODUCTION	4
2	ARCHITECTURE	5
2.1	GENERAL ARCHITECTURE	5
2.2	CODE ARCHITECTURE	6

1 Introduction

Next generation ***A***nalysis ***T***oolbox (NEAT) is a perl/R package that supports users during the analysis of next generation sequencing (NGS)

NEAT is versatile and easy to modify. In this tutorial, we will show how to add a custom module to NEAT. Adding a new module has been made very easy by automating all the repetitive tasks such as job creation, batch submission and queuing. Adding a new module usually falls down to a single line of code.

2 Architecture

2.1 General architecture

NEAT contains different modules (yellow boxes) than can be modified / added.

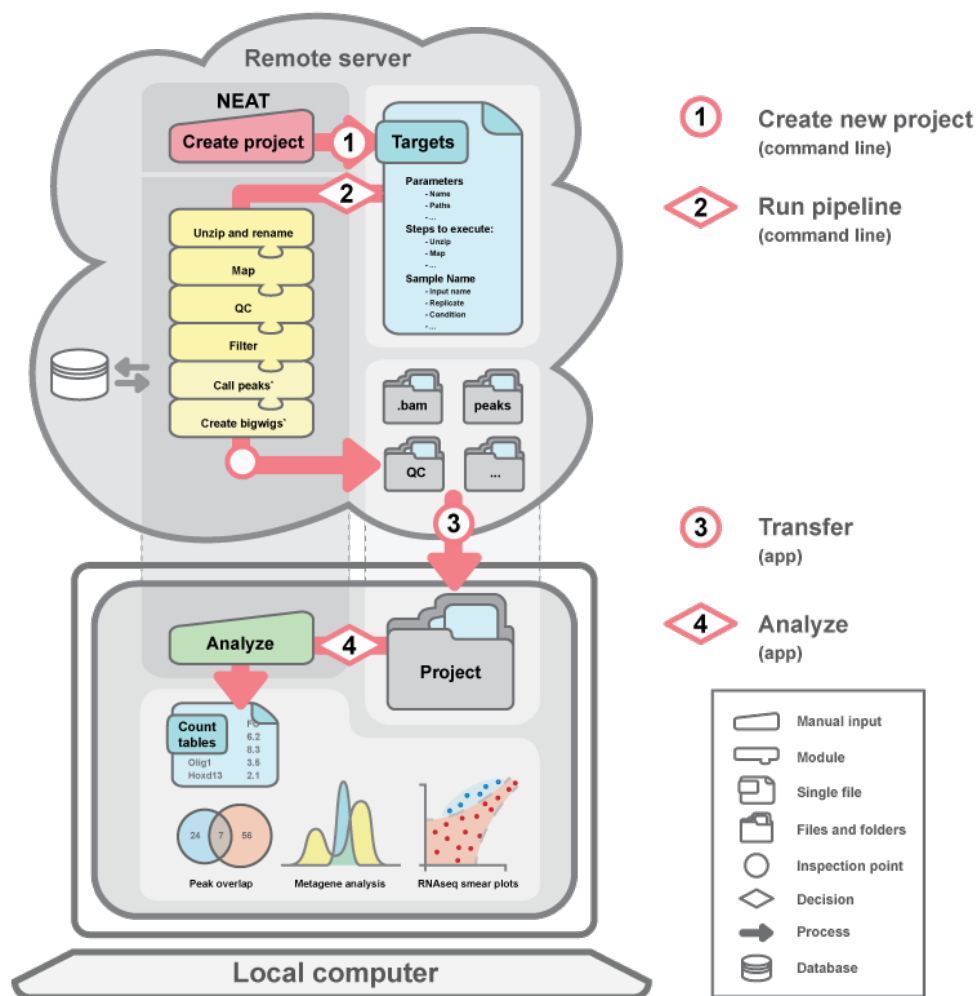


Fig.1 NEAT architecture. NGS data can be analyzed using NEAT in less than a day. Users follow a logical 4-step process, including the creation of a new project, running the pipeline on a remote server or in the cloud, transferring the data to a local computer and proceeding to the analysis.

2.2 Code architecture

The main code of NEAT (in the example of ChIPseq projects) is found in *./NEAT/ChIPpip/scripts/ChIPpip.pl*. The code is well annotated, highly redundant and should be self-explanatory to advanced users. Main modules are easily identifiable and customizable. A brief summary of how each module is built is depicted below.


```
# Mapping sequences with base
# Setup ~~~ "TRIP" ~~~
print("in Mapping func(job)"):
    my StrychJobName = "Iterate.<YOUR_JOB_NAME_HERE>"
    my StrychName = "<YOUR_JOB_NAME_HERE>"
    my StrychJobpath = "Strych/StrychJobName/Jobpath"

    # Create file to store job in
    unless [-d "Strych/StrychJobName"] { mkdir Strych/StrychJobName }
    unless [-d "StrychJobpath"] { mkdir StrychJobpath }
    my $SQUID = "Strych/StrychJobName/StrychJobName.squid"
    open $SQUID, ">" || "SQUID" or die "Can't open '$SQUID'"
    print $SQUID "in/job/jobpath\n"
    close $SQUID
    chmod 777 $SQUID
    print "I've Store all of the following 'StrychJobName' jobs in $SQUID " "n"
    my @myJobs

foreach my $ { 0..$StreamlinesInputs } {
    # Prepare a personal squid script
    my $SQUIDDir = "Strych/StrychJobName/StreamlinesInputs"
    mkdir $SQUIDDir

    # Create a directory
    my $pathToCurrentSampleDir = "StrychJobpath/$StreamlinesInputs"
    mkdir $pathToCurrentSampleDir

    #~~~~~
    #~~~~~
    my $cmd = "<YOUR_CUSTOM_COMMAND>"
    echo "$cmd" >> $SQUIDDir

    #~~~~~
    #~~~~~
    # Keep track of the jobs in @myJobs
    my $JobName = "StrychJobName/$"
    push(@myJobs, $JobName)
    $cmd = "StrychJobName/jobpath/$pathToCurrentSampleDir/$"
    open $SQUID, ">" || "SQUID" or die "Can't open '$SQUID'"
    print $SQUID "$cmd\n"
    close $SQUID
}

#~~~~~
# Change Targets.txt file for next iteration
print "in Changing 'StrychJobName' variable to FALSE and proceed"
"/usr/bin/perl -p -e s/'$StrychJobName'/StrychJobName_FALSE/g" "Targets.txt"

# Prepare file containing the jobs to run
# Add the next job line to the StrychJobDir
foreach (@myJobs){ $i = "1".."5" }
my $StrychJobDir = "job" . $i . @myJobs
my $StrychCmd = "FINAL" . @myJobs . StrychJobName . StrychJobpath . StrychJobName
open $SQUID, ">" || "SQUID" or die "Can't open '$SQUID'"
print $SQUID "$StrychCmd\n"
close $SQUID

#~~~~~
# Submit jobs to run

print "in Submitting jobs to cluster:" "n"
my $SQUIDDir

# Bat script

print "in Exiting StrychJobName section with no known error " "n"
exit 0;
```

Submit job to cluster and exit