

A Proposal for Beast 2.0

Vision

Basic design

XML

Applications

Summary

Remco R. Bouckaert

remco@cs.{auckland|waikato}.ac.nz

Department of Computer Science

University of Auckland & University of Waikato

To provide tools for computational science that are

- *easy to use*, that is, well documented, have intuitive user interfaces with small learning curve.
- *open access*, that is, open source, open xml format, facilitating reproducibility of results, runs on many platforms.
- *easy to extend*, by having extensibility in design.

Scope

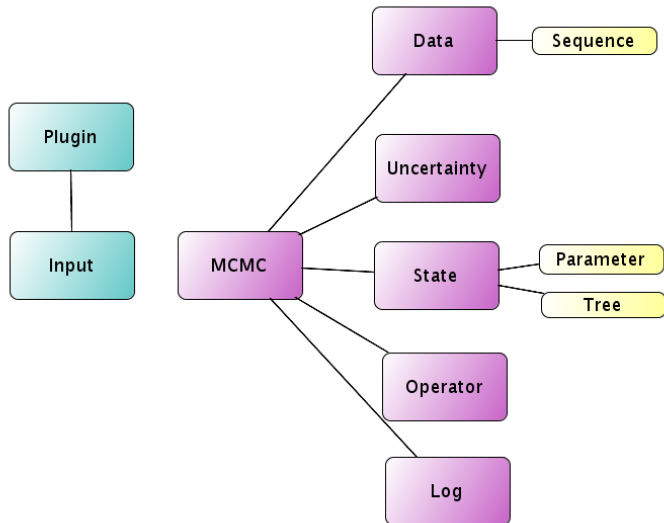
Efficient testing of probabilistic hypotheses for sequence data analysis involving tree models.

To provide tools for computational science that are

- *easy to use*, that is, **well documented**, have intuitive user interfaces with small learning curve.
- *open access*, that is, open source, open xml format, facilitating reproducibility of results, runs on many platforms.
- ***easy to extend***, by having extensibility in design.

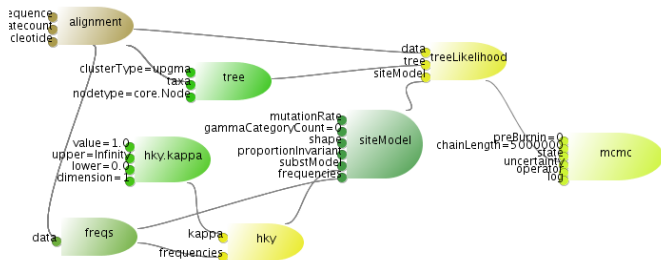
Scope

Efficient testing of probabilistic hypotheses for sequence data analysis involving tree models.



Phylosophy

Everything is a plug-in



Plug-ins provide...

- connection with other plug-ins/values through 'inputs'
- validation
- documentation
- 'XML parsing'

Plugin class

```
@Description("Description_goes_here")
public class Plugin {
    public void initAndValidate(State state)

    public String getCitations()

    public String getID()
    public void setID(String sID)

    public void store(int nSample)
    public void restore(int nSample)
} // class Plugin
```

HKY Plugin

```
@Description("HKY85_(Hasegawa,_Kishino_&_Yano,_1985)_substitution")
public class HKY extends 'Plugin' {
    public Input<Frequencies> m_freqs = new Input<Frequencies>("frequencies")
    public Input<Parameter> m_kappa = new Input<Parameter>("kappa")

    @Override public void initAndValidate(State state) throws Exception {
        initialiseEigen();
    }

    public void getTransitionProbabilities(double distance, double
        ...
    }
    @Override public void store(int nSample) {}
    @Override public void restore(int nSample) {
        updateMatrix = true;
        updateIntermediates = true;
    }
    @Override public String getCitation() {
        return "Hasegawa,_M.,_Kishino,_H_&_Yano,_T._1985._Dating
    }
} // class HKY
```

Beast 2.0

Bouckaert



Vision

Basic design

XML

throws Exception
Applications

Summary

Input validation

Default: OPTIONAL (see previous slide)

If input is REQUIRED:

```
public Input<Parameter> m_kappa =  
    new Input<Parameter>("kappa",  
        "kappa_parameter_in_HKY_model",  
        Validate.REQUIRED);
```

```
public Input<List<Operator>> m_operators =  
    new Input<List<Operator>>("operator",  
        "operator_for_generating_proposals_in_MCMC_state_space",  
        new ArrayList<Operator>(), Validate.REQUIRED);
```

If input is XOR:

```
public Input<Tree> m_pTree =  
    new Input<Tree>("tree",  
        "if_specified,_all_tree_branch_length_are_scaled");  
public Input<Parameter> m_pParameter =  
    new Input<Parameter>("parameter",  
        "if_specified,_this_parameter_is_scaled"  
        , Validate.XOR, m_pTree);
```

- State is explicit in XML & as object (unlike Beast 1)
- Contains parameters and trees
- Operators work on the state

```
public double proposal(State state) throws Exception { .. }
```

- Uncertainty calculated given state

```
public double calculateLogP(State state) throws Exception { .. }
```

- State can be interrogated on value of parameters/trees

```
state.getParameter(/**Parameter**/ p)  
state.getTree(/**Tree**/ t)
```

BEAST 1: push model

- Model
- Event handling
- Framework and plug-ins responsible
- Elegant, undocumented, requires lots of Panadol

BEAST 2: pull model

- store/restore
- sample number provided to prevent multiple (re)stores
- MCMC only tells its 'uncertainties' to store/restore, plug-ins responsible

```
<mcmc id="mcmc" chainLength="500000">
  <state>
    <parameter id="hky.kappa" value="1.0" lower="0.0"/>
    <tree spec='util.ClusterTree' id='tree' clusterType='upgma'>
      <input name='taxa' idref='alignment'/>
    </tree>
  </state>

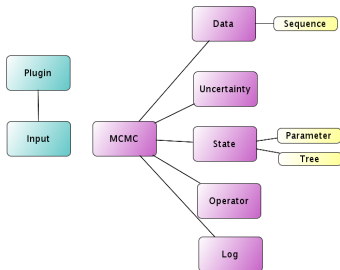
  <uncertainty id='likelihood' idref="treeLikelihood"/>

  <operator spec='ScaleOperator' scaleFactor="0.5" weight="1">
    <parameter idref="hky.kappa"/>
  </operator>
  <operator spec='ScaleOperator' scaleFactor="0.5" weight="1">
    <tree idref="tree"/>
  </operator>
  <operator spec='Uniform' weight="10">
    <tree idref="tree"/>
  </operator>
  <operator spec='SubtreeSlide' weight="5" gaussian="true" size="1.0">
    <tree idref="tree"/>
  </operator>
  <operator id='narrow' spec='Exchange' isNarrow='true' weight="1">
    <tree idref="tree"/>
  </operator>
  <operator id='wide' spec='Exchange' isNarrow='false' weight="1">
```

XML parsing/writing provided by the framework XML

Reserved elements

<mcmc >
<uncertainty >
<operator >
<log >
<data >
<sequence >
<state >
<parameter >
<tree >



```

<beast version='2.0' namespace='x.y.z:' >
<map name='elementName' >x.y.z.Class </map>
  
```

[Vision](#)
[Basic design](#)
[XML](#)
[Applications](#)
[Summary](#)

```
<!-- The HKY substitution model (Hasegawa, Kishino & Yano, 1985) -->
<input spec='HKY' id="hky">
  <parameter name='kappa' idref='hky.kappa' />
  <input id='freqs' name='frequencies' spec='Frequencies'>
    <input name='data' idref='alignment' />
  </input>
</input>

<!-- site model -->
<input spec='SiteModel' id="siteModel">
  <input idref='freqs' name='frequencies' />
  <input name='substModel' idref='hky' />
</input>

<input spec='TreeLikelihood' id="treeLikelihood">
  <input name='data' idref="alignment" />
  <input name='tree' idref="tree" />
  <input name='siteModel' idref="siteModel" />
</input>
```

```
<!-- The sequence alignment (each sequence refers to a taxon above). -->
<!-- ntax=6 nchar=768 -->
<!-- npatterns=69 -->
<data id="alignment" dataType="nucleotide">
  <sequence taxon="human">
AGAAAT ATGT CT GAT AAAAGAGTT ACTTT GAT AGAGT AAAT AAT AGGAGCTT AAACCCCCTT ATTT CT ACT AGGACT AT GAGAAT CGAAC
  </sequence>
  <sequence taxon="chimp">
AGAAAT ATGT CT GAT AAAAGAATT ACTTT GAT AGAGT AAAT AAT AGGAGTT CAAAT CCCCTT ATTT CT ACT AGGACT AT AAGAAT CGAAC
  </sequence>
  <sequence taxon="bonobo">
AGAAAT ATGT CT GAT AAAAGAATT ACTTT GAT AGAGT AAAT AAT AGGAGTTT AAAT CCCCTT ATTT CT ACT AGGACT AT GAGAGT CGAAC
  </sequence>
  <sequence taxon="gorilla">
AGAAAT ATGT CT GAT AAAAGAGTT ACTTT GAT AGAGT AAAT AAT AGAGGTTT AAACCCCCTT ATTT CT ACT AGGACT AT GAGAATT GAAC
  </sequence>
  <sequence taxon="orangutan">
AGAAAT ATGT CT GACAAAAGAGTT ACTTT GAT AGAGT AAAAAAT AGAGGT CT AAAT CCCCTT ATTT CT ACT AGGACT AT GGAATT GAAC
  </sequence>
  <sequence taxon="siamang">
AGAAAT ACGT CT GACGAAAGAGTT ACTTT GAT AGAGT AAAT AACAGGGGTTT AAAT CCCCTT ATTT CT ACT AGAACCAT AGGAGT CGAAC
  </sequence>
</data>
```

Reserved attributes:

```
<input id='myId'  
      idref='otherId'  
      name='inputName'  
      spec='x.y.z.MyClass' />
```

Resolving plug-in class:

- specified in spec attribute
- if not, get from element2class map
- if not, use element name (and hope it shows up in the namespace somewhere).

Resolving input name:

- specified in name attribute

```
<input name="xyz" >
```

- if not, use (non-reserved) attribute name

```
<input xyz="3" >
```

- if not, use element name

```
<xyz value="3" >
```

- if input, use 'value' when there is text content, but no element content

```
<input>3</input>
```

Resolving input value:

- if idref is specified, use the referred object

```
<input idref="other" >
```

- specified in value attribute

```
<xyz value="3" >
```

- if not, use value of (non-reserved) attribute

```
<input xyz="3" >
```

- if not, use text content when there is text content, but no element content

```
<input>3</input>
```

Parsing rules: Processing non reserved attributes

```
<input otherAttribute="xyz" />
```

equals

```
<input >  
  <input name='otherAttribute' value='xyz' />  
</input>
```

Processing non reserved element names

```
<myElement />
```

==

```
<input spec='myElement' name='myElement' />
```

unless 'spec' is a specified attribute, then that overrides,
likewise for 'name'

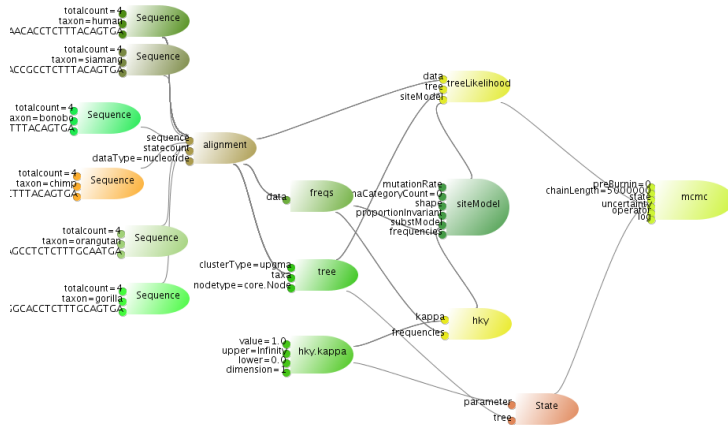
Processing of text content (only when there are no
enclosing elements)

```
<input name='data' >xyz </input>
```

==

```
<input name='data' value='xyz' / >
```

Model builder



Documentation

XML documentation provided through

- @Description annotation on plug in
- Tooltip text on inputs
- getCitation method
- Input validation rules

BEAST 2.0 Documentation - Mozilla Firefox

file:///home/r/b/workspace/GenerationD/doc/html/index.html

OM NOM NOM NOM


BEAST 2.0 Documentation index

- core
- [.Data](#)
- [.Logger](#)
- [.MCMC](#)
- [.Node](#)
- [.Operator](#)
- [.Parameter](#)
- [.Plugin](#)
- [.Sequence](#)
- [.State](#)
- [.Tree](#)
- [.Uncertainty](#)

- nuc
- [.CompoundUncertainty](#)
- [.Frequencies](#)
- nuc.likelihood
- [.TreeLikelihood](#)
- nuc.operators
- [.Exchange](#)
- [.ScaleOperator](#)
- [.SubtreeSlide](#)
- [.Uniform](#)
- [.WilsonBalding](#)
- nuc.sitemodel
- [.SiteModel](#)

BEAST 2.0 Documentation - Mozilla Firefox

file:///home/rb/workspace/GenerationD/doc/html/index.html



BEAST 2.0 Documentation: nuc.substitutionmodel.HKY

Specifies transition probability matrix for a given distance

HKY85 (Hasegawa, Kishino & Yano, 1985) substitution model of nucleotide evolution.

Reference:

Hasegawa, M., Kishino, H and Yano, T. 1985. Dating the human-ape splitting by a molecular clock of mitochondrial DNA. *Journal of Molecular Evolution* 22:160-174.

Inputs:

kappa

type: core.Parameter
kappa parameter in HKY model
Required input

frequencies

type: nuc.Frequencies
frequencies of characters
Required input

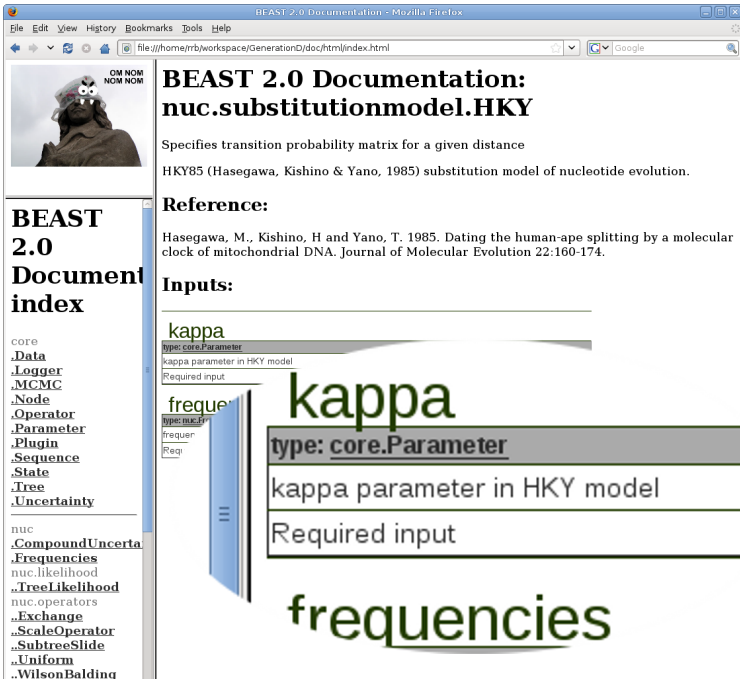
BEAST 2.0 Document index

core

[.Data](#)
[.Logger](#)
[.MCMC](#)
[.Node](#)
[.Operator](#)
[.Parameter](#)
[.Plugin](#)
[.Sequence](#)
[.State](#)
[.Tree](#)
[.Uncertainty](#)


nuc

[.CompoundUncertainty](#)
[.Frequencies](#)
[nuc.likelihood](#)
[..TreeLikelihood](#)
[nuc.operators](#)
[..Exchange](#)
[..ScaleOperator](#)
[..SubtreeSlide](#)
[..Uniform](#)
[..WilsonBalding](#)



BEAST 2.0 Documentation - Mozilla Firefox

file:///home/rb/workspace/GenerationD/doc/html/index.html



BEAST 2.0 Documentation: nuc.substitutionmodel.HKY

Specifies transition probability matrix for a given distance

HKY85 (Hasegawa, Kishino & Yano, 1985) substitution model of nucleotide evolution.

Reference:

Hasegawa, M., Kishino, H and Yano, T. 1985. Dating the human-ape splitting by a molecular clock of mitochondrial DNA. *Journal of Molecular Evolution* 22:160-174.

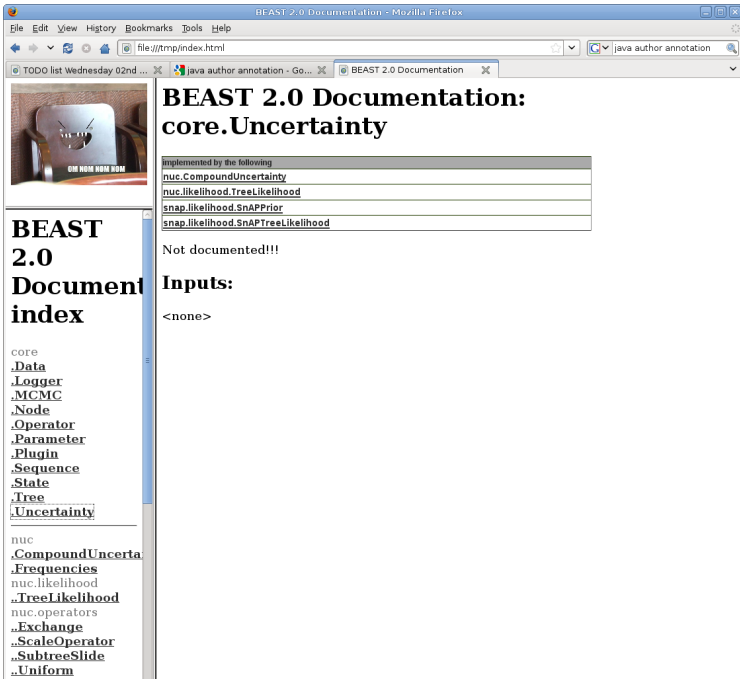
Inputs:

kappa
type: core.Parameter
kappa parameter in HKY model
Required input

frequencies
type: nuc.Frequencies
frequencies parameter in HKY model
Required input

BEAST 2.0 Document index

- core
 - .Data
 - .Logger
 - .MCMC
 - .Node
 - .Operator
 - .Parameter
 - .Plugin
 - .Sequence
 - .State
 - .Tree
 - .Uncertainty
- nuc
 - .CompoundUncertainty
 - .Frequencies
 - nuc.likelihood
 - ..TreeLikelihood
 - nuc.operators
 - ..Exchange
 - ..ScaleOperator
 - ..SubtreeSlide
 - ..Uniform
 - ..WilsonBalding




BEAST 2.0 Documentation - Mozilla Firefox

file:///tmp/index.html

java author annotation

TODD list Wednesday 02nd ... java author annotation - Go... BEAST 2.0 Documentation



BEAST 2.0 Documentation: core.Uncertainty

implemented by the following

nuc.CompoundUncertainty
nuc.likelihood.TreeLikelihood
snap.likelihood.SnAPPrior
snap.likelihood.SnAPTreeLikelihood

Not documented!!!

Inputs:

<none>

BEAST 2.0 Documentation index

- core
 - [.Data](#)
 - [.Logger](#)
 - [.MCMC](#)
 - [.Node](#)
 - [.Operator](#)
 - [.Parameter](#)
 - [.Plugin](#)
 - [.Sequence](#)
 - [.State](#)
 - [.Tree](#)
 - [.Uncertainty](#)
- nuc
 - [.CompoundUncertainty](#)
 - [.Frequencies](#)
 - [nuc.likelihood](#)
 - [.TreeLikelihood](#)
- nuc.operators
 - [.Exchange](#)
 - [.ScaleOperator](#)
 - [.SubtreeSlide](#)
 - [.Uniform](#)


```
~> java -cp bin app.BeastMCMC
```

Usage: BeastMCMC [options] <Beast.xml>

where <Beast.xml> the name of a file specifying a Beast run
and the following options are allowed:

-seed <int> : sets random number seed (**default** 127)

-threads <int> : sets number of threads (**default** 1)

BEAST 2.0

Beast 2.0

Bouckaert



[Vision](#)

[Basic design](#)

[XML](#)

[Applications](#)

[Summary](#)

HKY + nucleotide data, 'random' initial tree
2000 samples in debug mode, scaling if required, auto
optimise on, same logging, operators. etc.
testMCMC.xml – great apes

6 taxa

768 sites

69 patterns

10M samples single thread

	Beast 1.6/java	Beast1.6/native	Beast 2.0
real	3m55.056s	3m38.670s	2m31.794s
user	3m54.839s	3m38.670s	2m32.162s
sys	0m0.392s	0m0.428s	0m0.476s

BEAST 2.0

testMCMC.xml –

46 taxa
1363 sites
199 patterns
500K samples

Beast1.6/native core

real 1m56.097s
user 1m58.683s
sys 0m0.332s

Beast 2.0 + auto

real 0m56.843s
user 0m59.264s
sys 0m0.428s

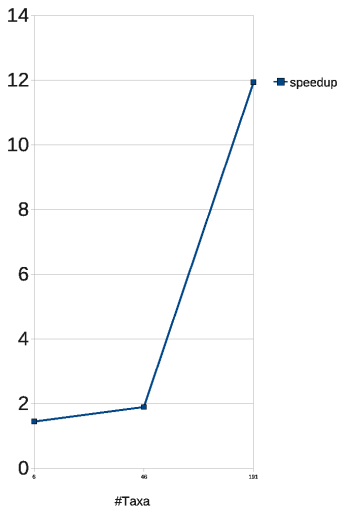
BEAST 2.0

testMCMC.xml –

191 taxa
606 sites
228 patterns
100K samples

Beast1.6/native
real 8m32.418s
user 8m33.060s
sys 0m0.476s

Beast 2.0
real 0m43.962s
user 0m46.051s
sys 0m0.432s



Summary Beast 2.0

Extensibility

- Everything is a plug-in
- Plug-in implementation:
 - Add @Description annotation
 - Specify Inputs
 - Implement initAndValidate
 - Optional: getCitation, store, restore
- Get xml-parsing, Input validation, xml-documentation for free
- Store/restore more transparent
- Explicit state object

Documentation

- Framework 'forces' proper documentation habits
- Automatic documentation generation

Performance

- Up to 1 order of magnitude better performance