atgtcttttctacaaaattttcatata agtcccggccaaacaataagatatgcagctaataatgcaggctctgtat aatactattaacqqcaataataataat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgtal aatactattaacggcaataataataa aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataa aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgta aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgl



# Recent developments in the Sequence Ontology

Colin Batchelor Royal Society of Chemistry, UK <u>batchelorc@rsc.org</u>

> Karen Eilbeck University of Utah, US eilbeck@fruitfly.org

atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataa aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaaccg atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgta aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatata atgcagctaataatgcaggctctgt

## The Sequence Ontology Project

### Overview

- What is SO?
- More-biologically-motivated relations
- Alignment with RO and topological relations
- Alignment with BFO: ontology curation experiments
- Open discussion

atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataa aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgta aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatata atgcagctaataatgcaggctctgta

# The Sequence Ontology Project

## What is SO? (1)

The Sequence Ontology organizes the kinds of, parts of and properties of biological sequence.

It deals with universal types.

The instances may be DNA molecules, RNA molecules or polypeptides. atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacqqcaataataataat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatataatgcagctaataatgcaggctctgt

# The Sequence Ontology Project

## What is SO? (2)

SO has been a pioneer in computable (cross-product) definitions of terms.

```
[Term]
id: S0:0000078
name: polycistronic_transcript
def: "A transcript that is polycistronic." [S0:xp]
intersection_of: S0:0000673 ! transcript
intersection_of: has_quality S0:0000880 ! Polycistronic
```

We are not restricted to dependent continuants for our differentiae:

```
[Term]
id: S0:0000111
name: transposable_element_gene
def: "A gene encoded within a transposable element. For example gag, int,
        env and pol are the transposable element genes of the TY element in
        yeast." [S0:ke]
intersection_of: S0:0000704 ! gene
intersection_of: part_of S0:0000101 ! transposable_element
```

atgtcttttctacaaaattttcatata agtcccggccaaacaataagatatgcagctaataatgcaggctctgtat aatactattaacqqcaataataataat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagat tatagagtttcaagtggaatactgcca aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacc agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgtal aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctg

# The Sequence Ontology Project

## Relations: the original set



is\_a:OK

part\_of: probably
not RO compatible in
some cases
(some-all instead of
all-some)

derives\_from: not
well-defined

atgtcttttctacaaaattttcatata agtcccggccaaacaataagatatgcagctaataatgcaggctctgtat aatactattaacqqcaataataataat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagatatgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacc tcgcatcagaagccgtttctacaacco atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcc atgcagctaataatgcaggctctgtal aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctcte

# The Sequence Ontology Project

### Relations: a second attempt



atgtcttttctacaaaattttcatata agtcccggccaaacaataagatatgcagctaataatgcaggctctgtat aatactattaacqqcaataataataat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgtat tttacaaacaatccagtaaacggatat aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca aatactattaacggcaataataataat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgta

## The Sequence Ontology Project

### Relations: work in progress

- Relations

- - is\_a
- ----- overlaps
- ---- part\_of
- processed
- ----- regulates
- ----- similar\_to

- union of — variant of

	transitive	symmetric	reflexive	antisymmetric
disjoint	-	+	-	-
adjacent to	-	-	+	-
equal	+	-	+	+
inside	+	-	-	-
contains	+	-	-	-
covers	+	-	-	-
covered by	+	-	-	-
overlaps	-	-	-	-

```
5. contains
```

|--A----| |-B--|

Feature A contains feature B if A and B share interior sequence b of A's boundary coincides with B's interior.

```
6. covers
|--A--|
|--B----|
Feature A covers feature B if both share a common boundary and in
```

See working\_draft.obo on CVS: http://song.cvs.sourceforge.net/song/ontology atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataa aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgtal aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatataatgcagctaataatgcaggctctgl



## **Relations: homology**

name: directly\_descends\_from
def: "F directly\_descends\_from F" iff there
 are O, O" such that O and O" are organisms, F
 and F" are features, O" is a parent of O and F
 has been copied from F"."

#### name: homologous\_to

def: "A feature F is homologous to another feature
 F' if F descends\_from F" and F'
 descends from F"."

atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacqqcaataataataat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgta

# The Sequence Ontology Project

## Alignment with RO

We are making sure all relations operate in the all-some direction. This involves, for example, replacing <code>part\_of with has\_part where appropriate</code>.

```
[Term]
id: S0:0001250
name: fingerprint_map
def: "A fingerprint_map is a physical map composed of
   restriction fragments." [S0:ke]
[...]
is_a: S0:0001249 ! fragment_assembly
relationship: has part S0:0000412 ! restriction fragment
```

atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco caaaattttcatataagtcccggccaa agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataa aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcc atgcagctaataatgcaggctctgtal aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatata atgcagctaataatgcaggctctgta

# The Sequence Ontology Project

## Alignment with BFO (1)

Currently: sequence \_feature terms in SO are related to sequence \_attribute terms in SO by the has quality relation.

This is not compatible with BFO!

#### Aim: to classify the

sequence attribute terms in SO according to the classes in BFO.



atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat gccaatactattaacaatagagcagat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco caaaattttcatataagtcccggccaa agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataa aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgta aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatata atgcagctaataatgcaggctctgt

The Sequence Ontology Project

## Alignment with BFO (2)

Iterative procedure:

- Two annotators
- Select 30 attributes randomly
- Assign each term into a category (quality, disposition, role, function) without consulting the other annotator
- Add justification
- Calculate agreement adjusted for chance (the kappa statistic)
- Discuss and modify guidelines
- Repeat

The Sequence Ontology Project

## Alignment with BFO (3)

Outcomes (after three rounds):

atgtcttttctacaaaattttcatata agtcccggccaaacaataagattatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgtat cgcctactaaggcaactcccataac aatactattaacggcaataataat gccaatactattaacaatagagcagat tttacaaacaatccagtaaacggatat

aatgaaagcgaccatggaaggatgt

aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgtat cgcctactaaggcaactccccataac aatactattaacggcaataataataa

aatgaaagcgaccatggaaggatg

aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco

atgtcttttctacaaaattttcatataagtcccggccaaacaataagattatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgtal

aatactattaacggcaataataataa

aatgaaagcgaccatggaaggatgt

aatgtaacggactttaactatacacca

atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgta

aatactattaacggcaataataataat

aatgaaagcgaccatggaaggatg

aatgtaacggactttaactatacacca

caaaattttcatataagtcccggccaa

atgtcttttctacaaaattttcatata

atgcagctaataatgcaggctctg

- Tentative, sequence-specific definitions of BFO classes in working\_draft.obo.
- However, we still have poor agreement: 53% (κ = 0.42) on Q vs. D vs. R vs. F. 77% (κ = 0.46) on Q vs. realizable entity
- More work on definitions needed!

atgtcttttctacaaaattttcatataagtcccggccaaacaataagatatgcagctaataatgcaggctctgtat aatactattaacqqcaataataataat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatgt cacgcccctccgcctgaacaattacaa aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagatatgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaacco caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgtat gccaatactattaacaatagagcagat tttacaaacaatccagtaaacggata aatgaaagcgaccatggaaggatgt cacgcccctccgcctgaacaattacaa aatgtaacggactttaactatacacca tcgcatcagaagccgtttctacaaccc caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatataagtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgta

The Sequence Ontology Project

## Any questions?

atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat tttacaaacaatccagtaaacggatat aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatata agtcccggccaaacaataagat atgcagctaataatgcaggctctgtat aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacc tcgcatcagaagccgtttctacaacco atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctgta aatactattaacggcaataataataa tttacaaacaatccagtaaacggata aatgaaagcgaccatggaaggatgt aatgtaacggactttaactatacacca atgtcttttctacaaaattttcatata agtcccggccaaacaataagat tatagagtttcaagtggaatactgcca atgcagctaataatgcaggctctgtal aatactattaacggcaataataataat aatgaaagcgaccatggaaggatg aatgtaacggactttaactatacacca caaaattttcatataagtcccggccaa atgtcttttctacaaaattttcatataagtcccggccaaacaataagat atgcagctaataatgcaggctctg

# The Sequence Ontology Project

## The kappa statistic

## $c = \frac{\text{actual agreement} - \text{chance agreement}}{1 - \text{chance agreement}}$

Widely-used in computational linguistics, psychology, clinical medicine and social sciences to adjust for chance agreement between human annotators. Always lower than or equal to percentage agreement.

 $\kappa$  = 1 indicates perfect agreement.

 $\kappa$  = 0 indicates agreement no better than if people were selecting categories at random.

 $\kappa$  = 0.67 is a lower bound for "acceptable" agreement.