

**RNA**  
**SEQ**

THOUGHTS AND IDEAS TOWARD OPTIMAL



RNA's will works in  
mysterious ways

- PIPELINE
- ```

RNA-seq (100 bp - PE)
↓
1. CLEAN&P (?)
↓
2. NORMALIZE (?)
↓
3. ASSEMBLE
↓
4. REDUNDANCY
↓
5. ANNOTATE

```

| Assembly number | Trim quality (PHRED) | Normalized | Number of transcripts | Number of genes | Unigene NS6 |
|-----------------|----------------------|------------|-----------------------|-----------------|-------------|
| 1               | 20                   | Y          | 1 68 089              | 126 617         | 1 655       |
| 2               | 10                   | Y          | 1 68 499              | 126 565         | 1 631       |
| 3               | 5                    | Y          | 1 68 742              | 126 565         | 1 631       |
| 4               | 2                    | Y          | 1 68 797              | 129 759         | 1 693       |
| 5               | 20                   | N          | 1 73 373              | 145 460         | 1 425       |
| 6               | 10                   | N          | 1 73 373              | 145 460         | 1 425       |
| 7               | 5                    | N          | 1 75 326              | 145 410         | 1 426       |
| 8               | 2                    | N          | 1 76 279              | 145 407         | 1 427       |

PLUS

| Evaluation of de novo transcriptome assemblies from RNA-Seq data                                                                |  |  |  |  |  |
|---------------------------------------------------------------------------------------------------------------------------------|--|--|--|--|--|
| <a href="http://develab.biology.uvic.ca/deNovoAssemblerEvaluator/">http://develab.biology.uvic.ca/deNovoAssemblerEvaluator/</a> |  |  |  |  |  |
| Reference file: de_novo_transcriptome_assembly_evaluator                                                                        |  |  |  |  |  |

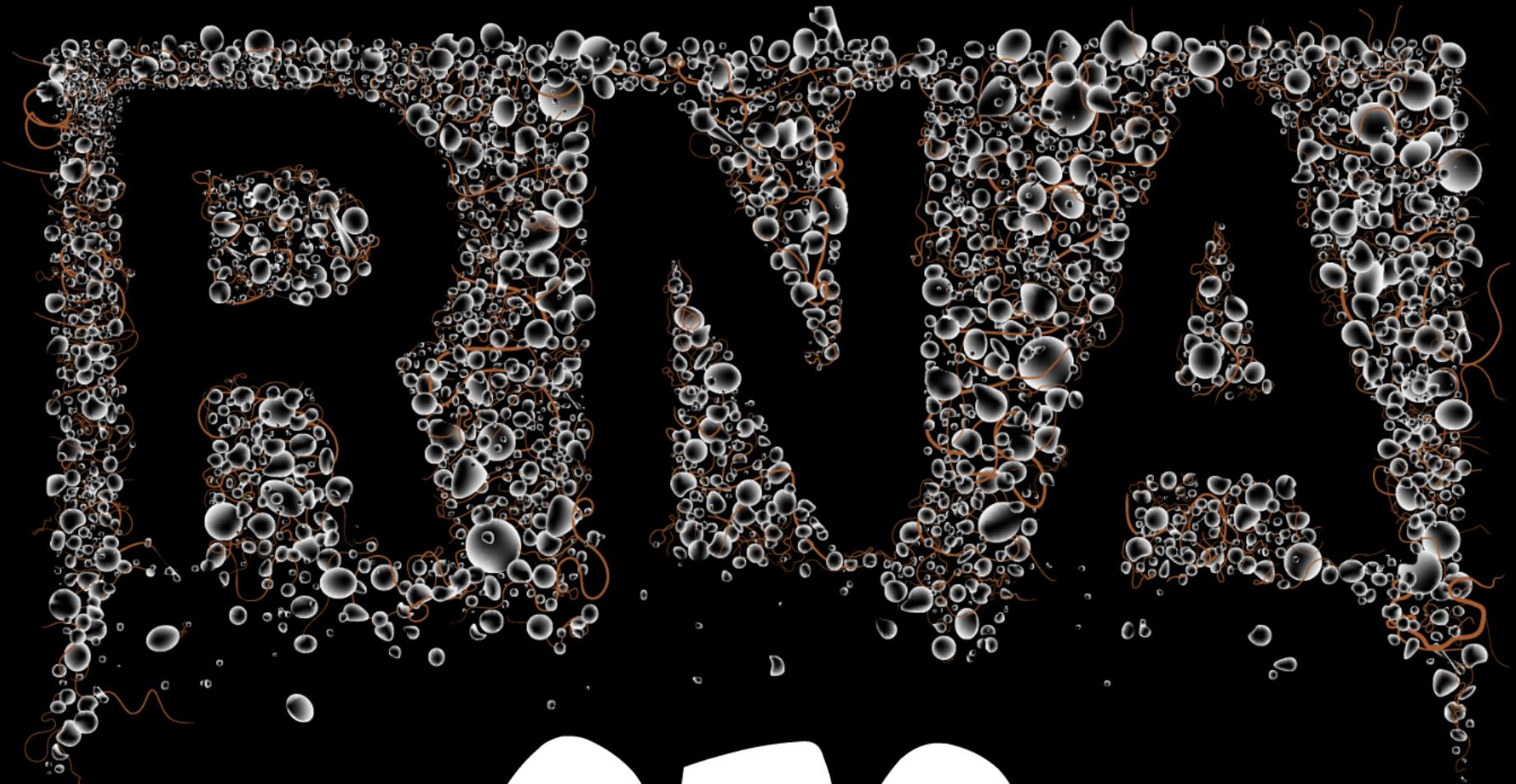
RESULTS SO FAR...

| Assembly number | Trim quality (PHRED) | Normalized | Number of transcripts | Number of genes | Unigene NS6 |
|-----------------|----------------------|------------|-----------------------|-----------------|-------------|
| 1               | 20                   | Y          | 1 68 089              | 126 617         | 1 655       |
| 2               | 10                   | Y          | 1 68 499              | 126 565         | 1 631       |
| 3               | 5                    | Y          | 1 68 742              | 126 565         | 1 631       |
| 4               | 2                    | Y          | 1 68 797              | 129 759         | 1 693       |
| 5               | 20                   | N          | 1 73 373              | 145 460         | 1 425       |
| 6               | 10                   | N          | 1 73 373              | 145 460         | 1 425       |
| 7               | 5                    | N          | 1 75 326              | 145 410         | 1 426       |
| 8               | 2                    | N          | 1 76 279              | 145 407         | 1 427       |

# reads = 19 370 403

# assembled bases = 4 682 256 300

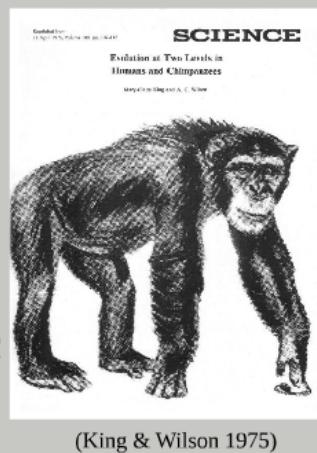
THOUGHTS AND IDEAS TOWARD OPTIMAL



SEQ

## RNA = powerful biological tool

- Response to the environment
- Ecological interactions
- Inter/Intra specific divergence

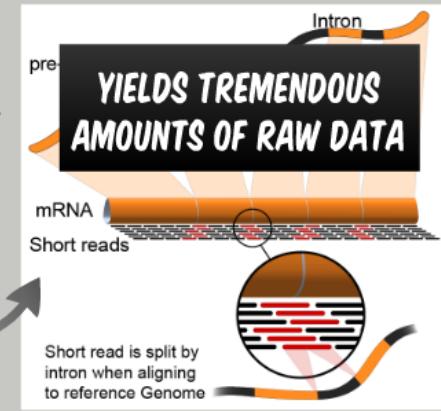


(King & Wilson 1975)



## GETTING RNA = EASY MONEY !

Any biological tissue OR whole organism



RNA

10% = extracting RNA

VS

90% = clean + annotate + DEBUG



## FRIENDLY ADVICE

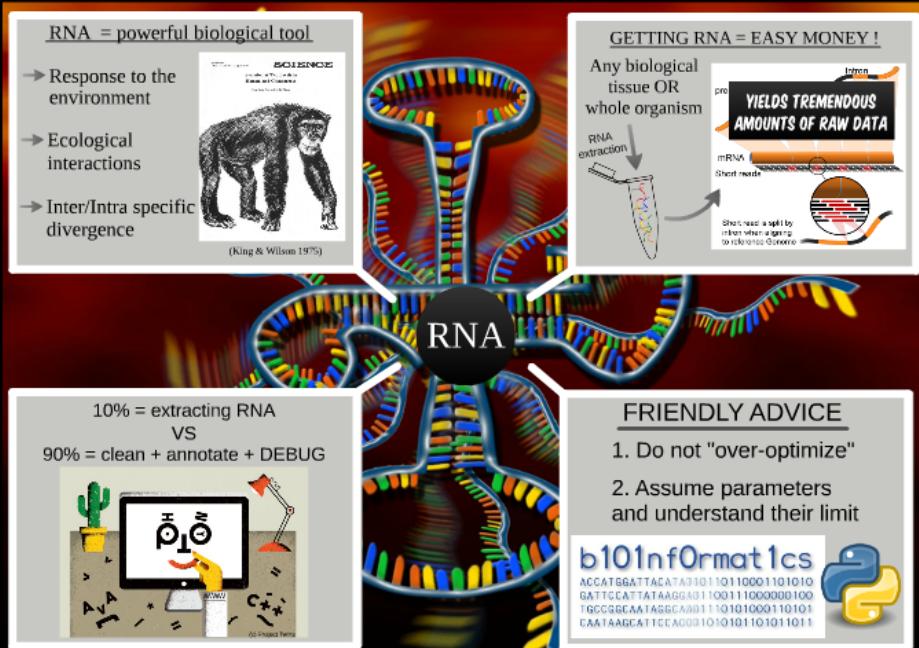
1. Do not "over-optimize"
2. Assume parameters and understand their limit

b101nf0rmat1cs

ACCATGGATTACATA@10110110001101010  
GATTC CATTATAAGGA01100111000000100  
TGCCGGCAATAGGC@01110101000110101  
CAATAAGCATTCCA0001010101101011011



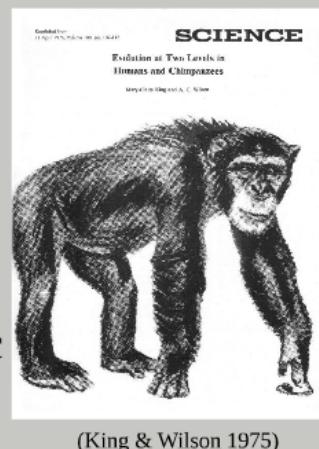
# RNA AS AN HINDU GODDESS ???



RNA's will works in  
mysterious ways

## RNA = powerful biological tool

- Response to the environment
- Ecological interactions
- Inter/Intra specific divergence



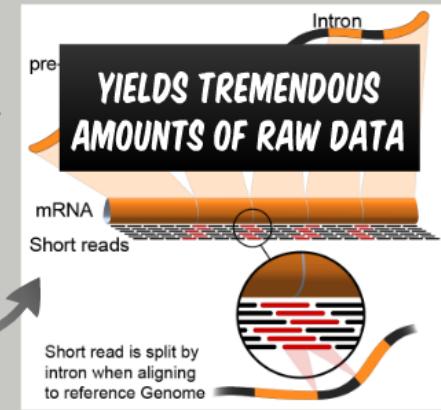
(King & Wilson 1975)



## GETTING RNA = EASY MONEY !

Any biological tissue OR whole organism

RNA extraction



## FRIENDLY ADVICE

1. Do not "over-optimize"
2. Assume parameters and understand their limit

b101nf0rmat1cs

ACCATGGATTACATA@10110110001101010  
GATTC CATTATAAGGA01100111000000100  
TGCCGGCAATAGGC@01110101000110101  
CAATAAGCATTCCA0001010101101011011

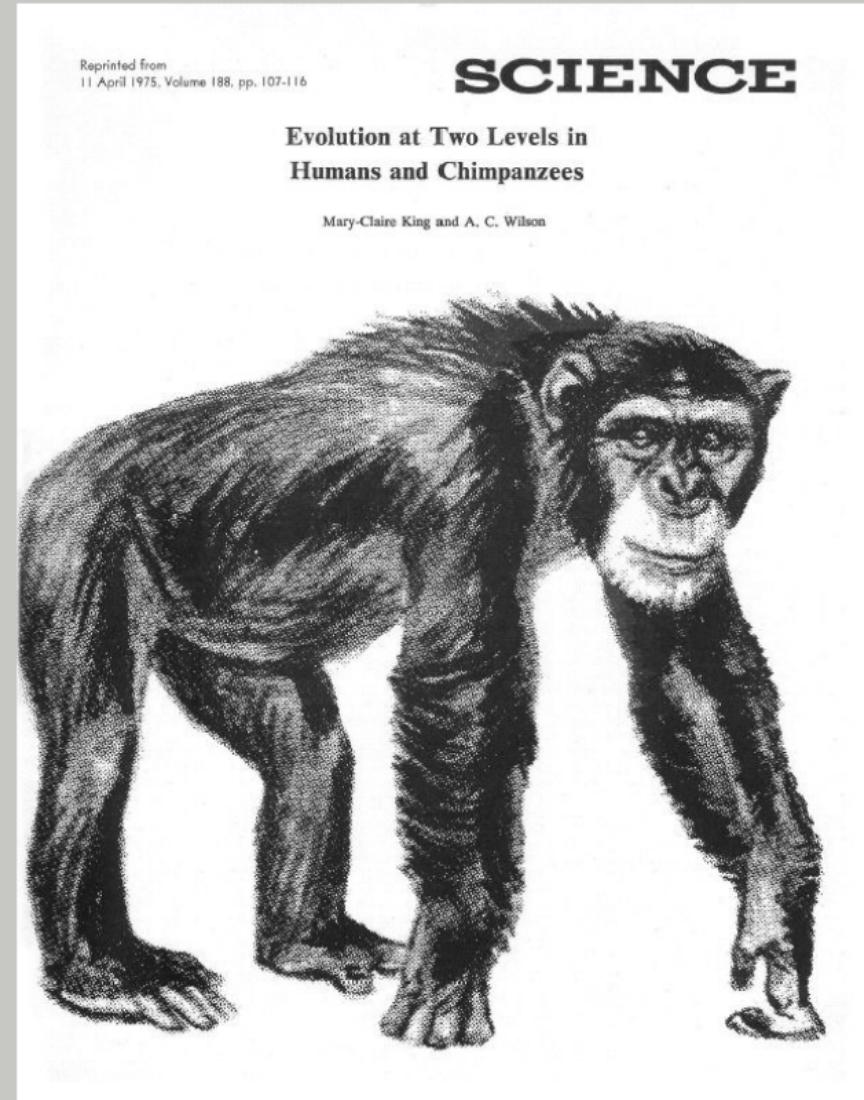


10% = extracting RNA  
VS  
90% = clean + annotate + DEBUG



# RNA = powerful biological tool

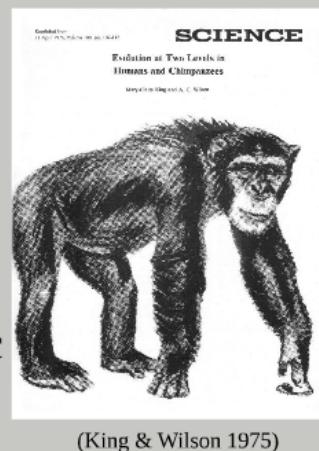
- Response to the environment
- Ecological interactions
- Inter/Intra specific divergence



(King & Wilson 1975)

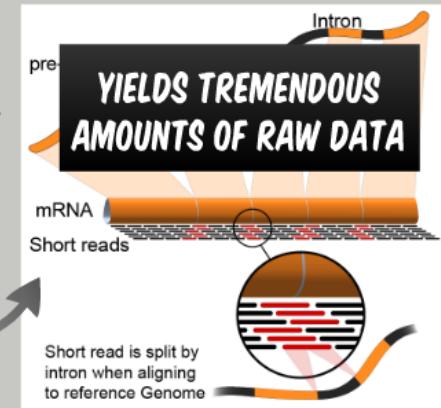
## RNA = powerful biological tool

- Response to the environment
- Ecological interactions
- Inter/Intra specific divergence



## GETTING RNA = EASY MONEY !

Any biological tissue OR whole organism



RNA

10% = extracting RNA

VS

90% = clean + annotate + DEBUG



## FRIENDLY ADVICE

1. Do not "over-optimize"
2. Assume parameters and understand their limit

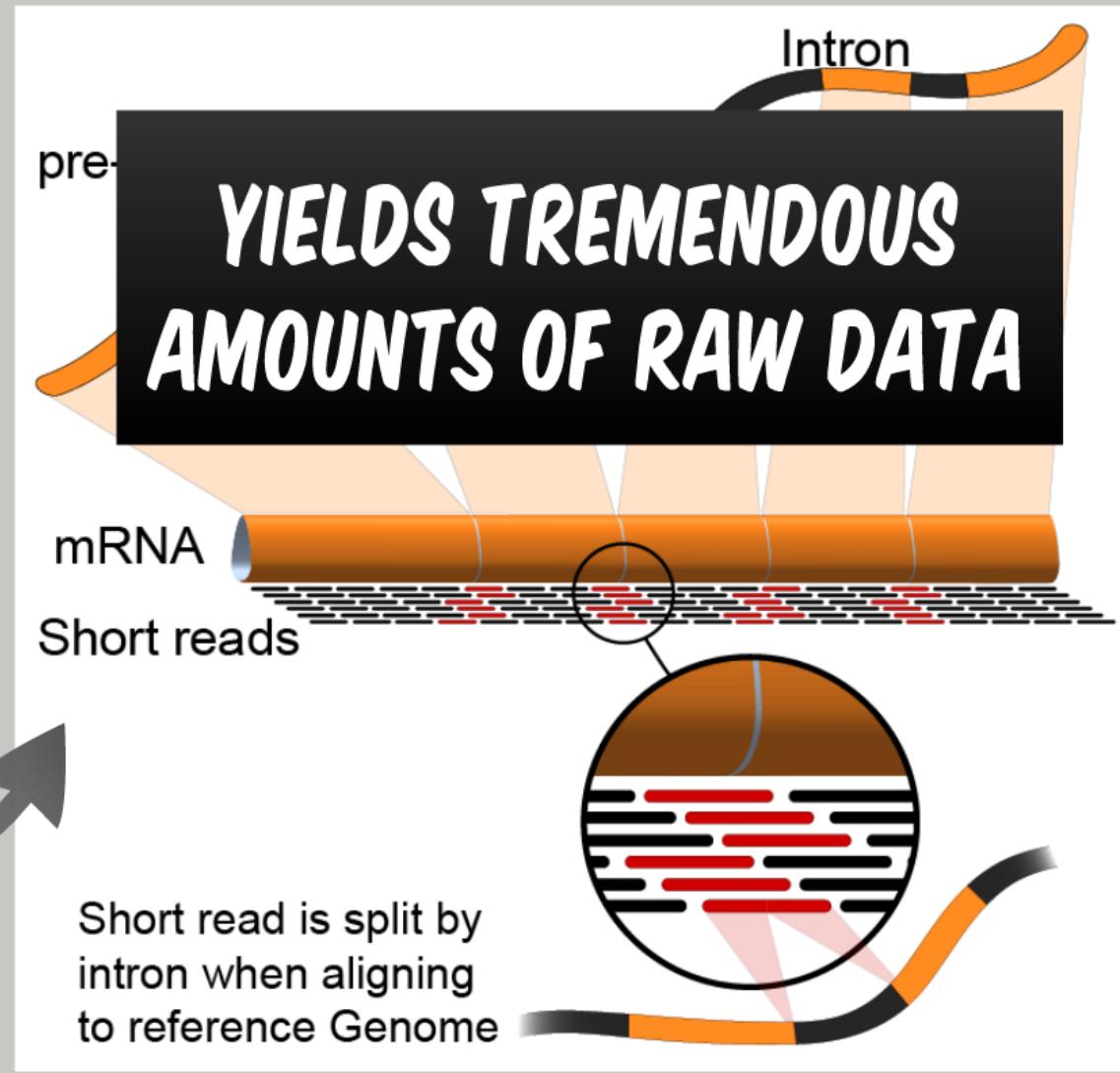
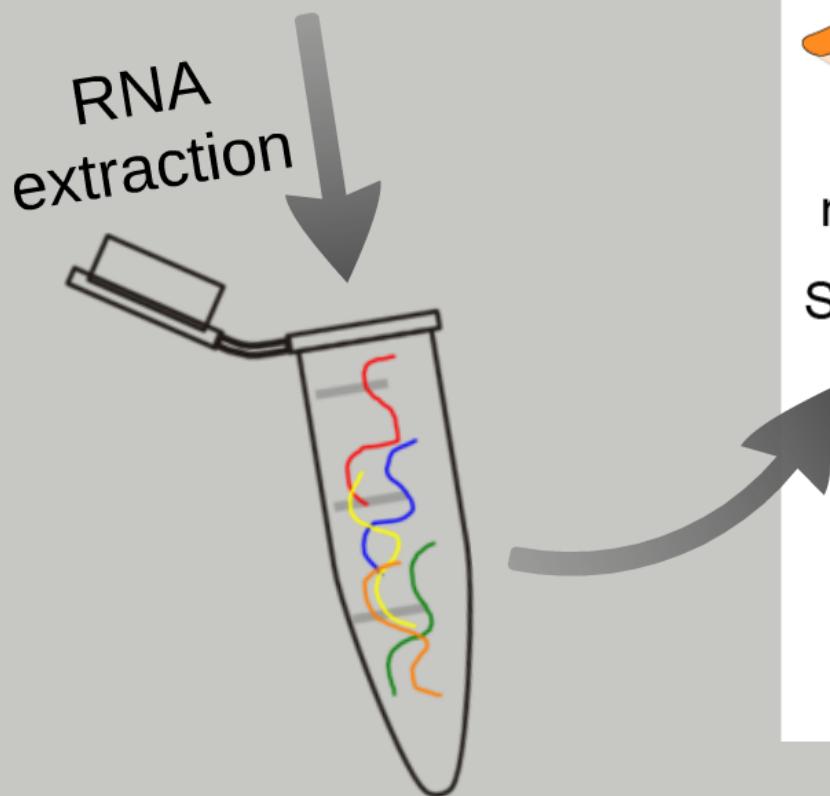
b101nf0rmat1cs

ACCATGGATTACATA@10110110001101010  
GATTC CATTATAAGGA01100111000000100  
TGCCGGCAATAGGC@01110101000110101  
CAATAAGCATTCCA0001010101101011011



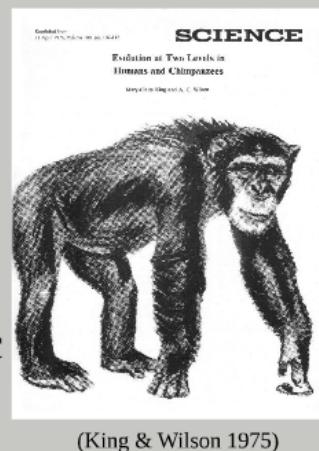
# GETTING RNA = EASY MONEY !

Any biological  
tissue OR  
whole organism



## RNA = powerful biological tool

- Response to the environment
- Ecological interactions
- Inter/Intra specific divergence



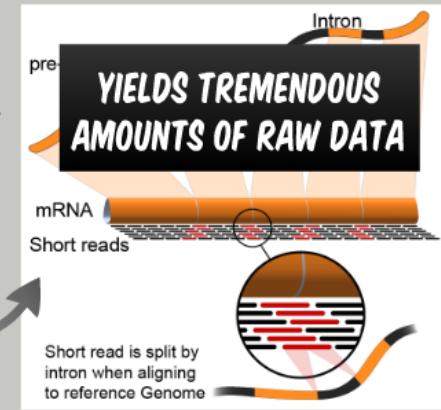
(King & Wilson 1975)



## GETTING RNA = EASY MONEY !

Any biological tissue OR whole organism

RNA extraction



RNA

10% = extracting RNA

VS

90% = clean + annotate + DEBUG



(c) Project Twins

## FRIENDLY ADVICE

1. Do not "over-optimize"
2. Assume parameters and understand their limit

b101nf0rmat1cs

ACCATGGATTACATA@10110110001101010  
GATTC CATTATAAGGA01100111000000100  
TGCCGGCAATAGGC@01110101000110101  
CAATAAGCATTCCA0001010101101011011

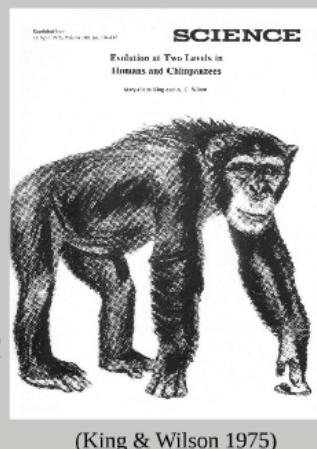


10% = extracting RNA  
VS  
90% = clean + annotate + DEBUG



## RNA = powerful biological tool

- Response to the environment
- Ecological interactions
- Inter/Intra specific divergence



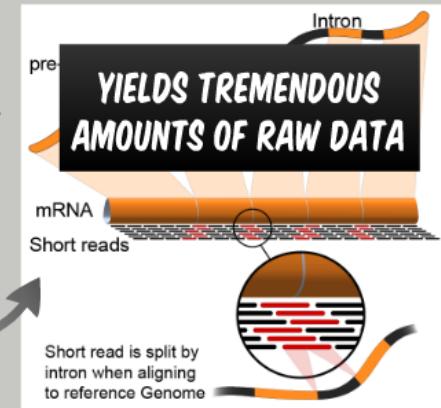
(King & Wilson 1975)



## GETTING RNA = EASY MONEY !

Any biological tissue OR whole organism

RNA extraction



## FRIENDLY ADVICE

1. Do not "over-optimize"
2. Assume parameters and understand their limit

b101nf0rmat1cs

ACCATGGATTACATA@10110110001101010  
GATTC CATTATAAGGA01100111000000100  
TGCCGGCAATAGGC@01110101000110101  
CAATAAGCATTCCA0001010101101011011



10% = extracting RNA  
VS  
90% = clean + annotate + DEBUG

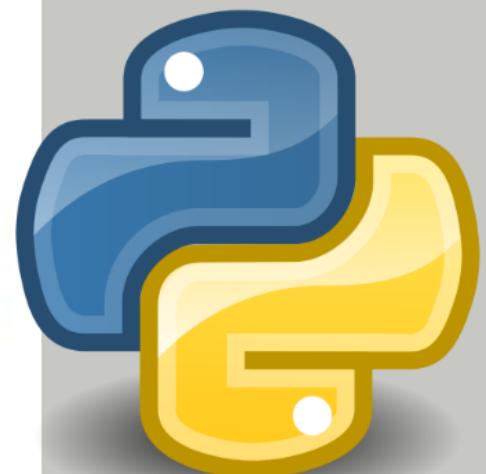


# FRIENDLY ADVICE

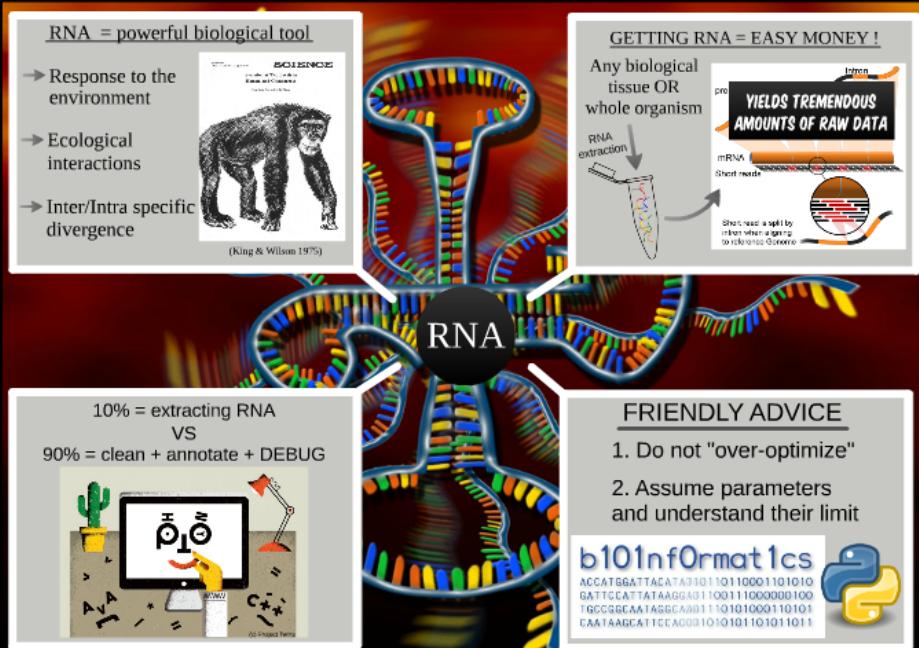
1. Do not "over-optimize"
2. Assume parameters  
and understand their limit

b101nf0rmat1cs

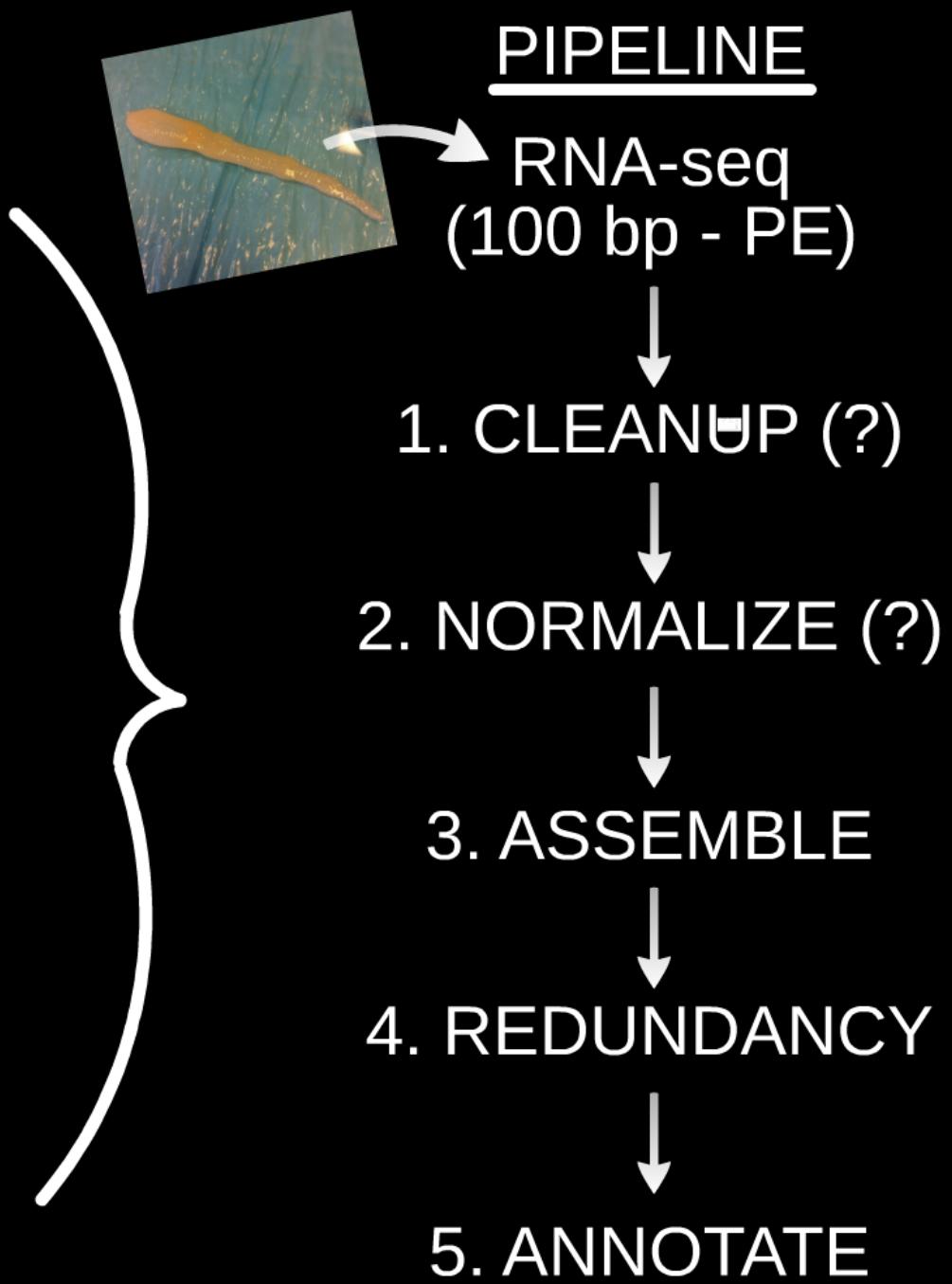
ACCATGGATTACATA@0110110001101010  
GATTCCATTATAAGGA@1100111000000100  
TGCGGGCAATAGGCA@01110101000110101  
CAATAAGCATTCCA@0001010101101011011



# RNA AS AN HINDU GODDESS ???



RNA's will work in  
mysterious ways



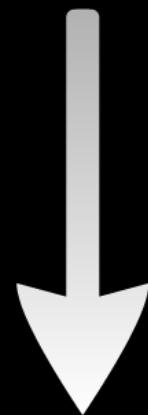


PIPELINE

RNA-seq  
(100 bp - PE)

1. CLEANUP (

RNA-seq  
(100 bp - PE)



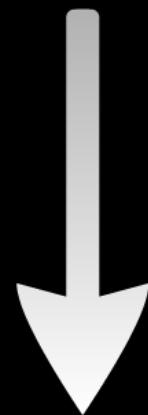
1. CLEANUP (?)



## Table 1: Quality Scores and Base Calling Accuracy

| Phred Quality Score | Probability of Incorrect Base Call | Base Call Accuracy |
|---------------------|------------------------------------|--------------------|
| 10                  | 1 in 10                            | 90%                |
| 20                  | 1 in 100                           | 99%                |
| 30                  | 1 in 1,000                         | 99.9%              |
| 40                  | 1 in 10,000                        | 99.99%             |
| 50                  | 1 in 100,000                       | 99.999%            |

RNA-seq  
(100 bp - PE)

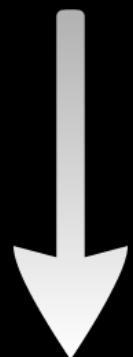


1. CLEANUP (?)





2. NORMALIZE (?)

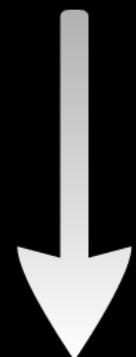


3. ASSEMBLE





4. REDUNDANCY



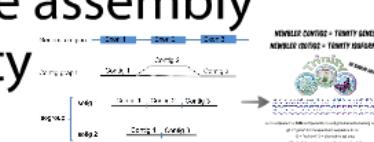
5. ANNOTATE

RESEARCH ARTICLE

Open Access

# Assessing *De Novo* transcriptome assembly metrics for consistency and utility

Shawn T O'Neil<sup>1,2</sup> and Scott J Emrich<sup>2\*</sup>



**Table 1 Metric trends and consistency**

| Metric                                 | Trend by sequencing depth | Consistent over sequencing depths | Trend by read length | Consistent over read lengths | Fully consistent metric |
|----------------------------------------|---------------------------|-----------------------------------|----------------------|------------------------------|-------------------------|
| Contig count                           | ↗                         | ✓                                 | ↘                    | ✓                            | ✗                       |
| % of reads used in contigs             | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| BP in contigs                          | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| % BP in contigs                        | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| Average contig coverage                | ↗                         | ✗                                 | ↗                    | ✗                            | ✗                       |
| Average unigene coverage               | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| Contig read count COV                  | ↘                         | ✓                                 | ↗                    | ✓                            | ✗                       |
| Unigene read count COV                 | ↘                         | ✓                                 | ↗                    | ✓                            | ✗                       |
| Average contig length                  | ↗                         | ✗                                 | ↗                    | ✗                            | ✗                       |
| Average unigene length                 | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| Contig N50 length                      | ↗                         | ✗                                 | ↗                    | ✗                            | ✗                       |
| Unigene N50 length                     | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| Unique annotations in singletons       | ↘                         | ✓                                 | ↗                    | ✓                            | ✗                       |
| Unique annotations in contigs          | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| Unique annotations in unigenes         | ↗                         | ✗                                 | ↗                    | ✓                            | ✗                       |
| Average contig OHR                     | ↗                         | ✗                                 | ↗                    | ✗                            | ✗                       |
| Average unigene OHR                    | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| Contig RBH count                       | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| Unigene RBH count                      | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| % of Annotated contigs with RBHs       | ↘                         | ✗                                 | ↗                    | ✗                            | ✗                       |
| % of Annotated unigenes with RBHs      | ↘                         | ✓                                 | ↗                    | ✗                            | ✗                       |
| Average contig CF                      | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| Average unigene CF                     | ↗                         | ✗                                 | ↗                    | ✗                            | ✗                       |
| Unique reverse annotations in contigs  | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |
| Unique reverse annotations in unigenes | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                       |

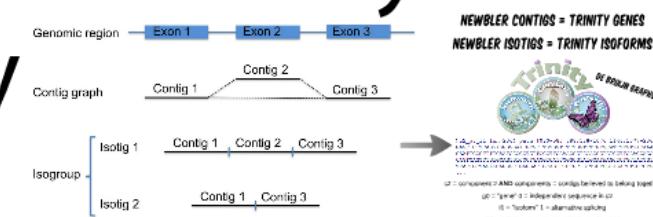
PLUS

RESEARCH ARTICLE

Open Access

# Assessing *De Novo* transcriptome assembly metrics for consistency and utility

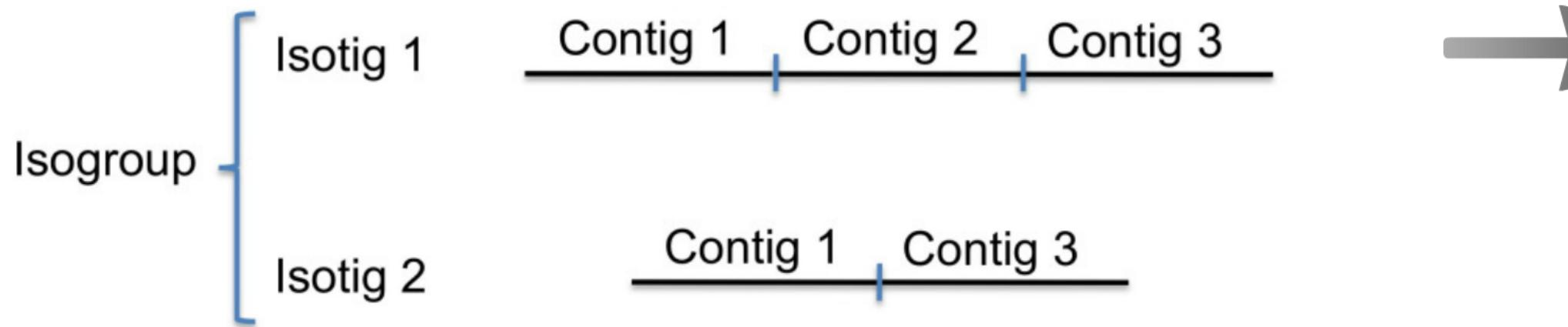
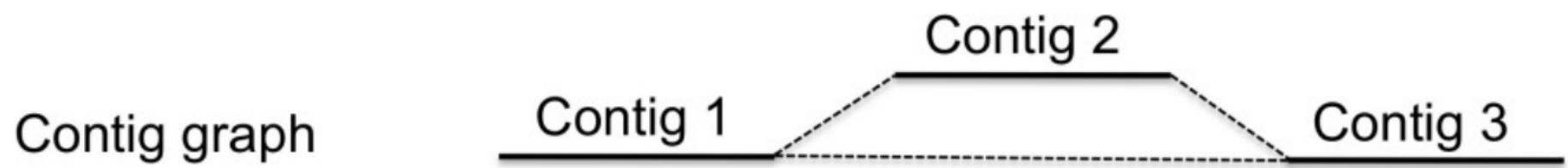
Shawn T O'Neil<sup>1,2</sup> and Scott J Emrich<sup>2\*</sup>



## Metric trends and consistency

| Metric                                                                                             | Trend by sequencing depth | Consistent over sequencing depths | Trend by read length | Consistent over read lengths | Fully consistent |
|----------------------------------------------------------------------------------------------------|---------------------------|-----------------------------------|----------------------|------------------------------|------------------|
| Contig count                                                                                       | ↗                         | ✓                                 | ↘                    | ✓                            | ✗                |
| of reads used in contigs = $\frac{\# \text{ BASES IN "GENES"} }{\# \text{ TOTAL ASSEMBLED BASES}}$ | ↗                         | ✓                                 | ↗                    | ✓                            | ✓                |

α β γ δ ε η ι ο ι γ





```
>c2_g0_i1 len=2364 path=[0:0-587 588:588-1076 1146:1077-2363]
GAGCTCTTCAGGAGGGGAATGTGCTTGTGGTTTTGGTCTTGTGCATTTGTGACAAAG
GAATTCCCTTTGAATCGCGCTGTTCCCTTGAAACCCCTGGAGCCTCTGGTTCAAGCAGCG
CAGTCAGTCTGTGCAGTGTCCCTGACGTCATCCGGCGTATGCATAAGCTCTGCTATTGTC
TTACCGCTAGAGCAGGGCTGAGGACTGCAGTCTGCTGCTCGCAGACCTGCCCTGC
...
```

c2 = component 2 **AND** components = contigs believed to belong together

g0 = "gene" 0 = independent sequence in c2

i1 = "isoform" 1 = alternative splicing

([https://github.com/enormandeau/trinity\\_pipeline\\_ibis](https://github.com/enormandeau/trinity_pipeline_ibis))

**NEWBLER CONTIGS = TRINITY GENES**

3

**NEWBLER ISOTIGS = TRINITY ISOFORMS**



*DE BRUIJN GRAPHS*



```
>c2_g0_i1 len=2364 path=[0:0-587 588:588-1076 1146:1077-2363]
GAGCTCTTCAGGAGGGGAATGTGCTTGTGGTTTGGTCTTGTGCATTTGTGACAAAG
GAATTCCCTTTGAATCGCGCTGTTCCCTTGAAACCCTGGAGCCTCTGGTTCAAGCAGCG
CAGTCAGTCTGTGCAGTGTCCCTGACGTCATCCGGCGTATGCATAAGCTCTGCTATTGTC
TTACCGCTAGAGCAGGGCTGAGGACTGCAGTCTCTGCTGCTGCTCGCAGACCTGCCCTGC
...
```

c2 = component 2 **AND** components = contigs believed to belong together

g0 = "gene" 0 = independent sequence in c2

i1 = "isoform" 1 = alternative splicing

([https://github.com/enormandeau/trinity\\_pipeline\\_ibis](https://github.com/enormandeau/trinity_pipeline_ibis))

**Table 1 Metric trends and consistency**

| Metric                                                                                                     | Trend by         | Consistent over   | Trend by    | Consistent over | Fully consistent |
|------------------------------------------------------------------------------------------------------------|------------------|-------------------|-------------|-----------------|------------------|
|                                                                                                            | sequencing depth | sequencing depths | read length | read lengths    | metric           |
| Contig count                                                                                               | ↗                | ✓                 | ↘           | ✓               | ✗                |
| % of reads used in contigs<br>= <small>* BASES IN "GENES"<br/>+ TOTAL ASSEMBLED BASES</small>              | ↗                | ✓                 | ↗           | ✓               | ✓                |
| BP in contigs                                                                                              | ↗                | ✓                 | ↗           | ✓               | ✓                |
| % BP in contigs                                                                                            | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Average contig coverage                                                                                    | ↗                | ✗                 | ↗           | ✗               | ✗                |
| Average unigene coverage                                                                                   | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Contig read count COV                                                                                      | ↔                | ✓                 | ↘           | ✓               | ✗                |
| Unigene read count COV                                                                                     | ↔                | ✓                 | ↘           | ✓               | ✗                |
| Average contig length                                                                                      | ↗                | ✗                 | ↗           | ✗               | ✗                |
| Average unigene length                                                                                     | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Contig N50 length                                                                                          | ↗                | ✗                 | ↗           | ✗               | ✗                |
| Unigene N50 length<br>= <small>* NUMBER OF READS OF LENGTH EQUAL TO OR LONGER THAN THE N50 LENGTH</small>  | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Unique annotations in singletons                                                                           | ↔                | ✓                 | ↘           | ✓               | ✗                |
| Unique annotations in contigs                                                                              | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Unique annotations in unigenes                                                                             | ↗                | ✗                 | ↔           | ✓               | ✗                |
| Average contig OHR                                                                                         | ↗                | ✗                 | ↗           | ✗               | ✗                |
| Average unigene OHR<br>= <small>* NUMBER OF READS OF LENGTH EQUAL TO OR LONGER THAN THE N50 LENGTH</small> | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Contig RBH count                                                                                           | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Unigene RBH count                                                                                          | ↗                | ✓                 | ↗           | ✓               | ✓                |
| % of Annotated contigs with RBHs                                                                           | ↔                | ✗                 | ↗           | ✗               | ✗                |
| % of Annotated unigenes with RBHs                                                                          | ↔                | ✓                 | ↗           | ✗               | ✗                |
| Average contig CF                                                                                          | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Average unigene CF                                                                                         | ↗                | ✗                 | ↗           | ✗               | ✗                |
| Unique reverse annotations in contigs                                                                      | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Unique reverse annotations in unigenes                                                                     | ↗                | ✓                 | ↗           | ✓               | ✓                |

---

# Contig count

---

% of reads used in contigs

---

# BP in contigs

---

% BP in contigs

---

**# BASES IN "GENES"**

=

---

**# TOTAL ASSEMBLED BASES**

---

---

# Contig count

---

% of reads used in contigs

---

# BP in contigs

---

% BP in contigs

Average contig coverage

---

Average unigene coverage

---

Contig read count COV

---

Unigene read count COV

---

# Average contig length

---

## Average unigene length

---

### Contig N50 length

---

#### Unigene N50 length

= LENGTH AT WHICH  
THIS LENGTH  
CONTAIN AT LEAST  
ASSEMBLED

ique annotations in singlet

= LENGTH AT WHICH CONTIGS OF  
THIS LENGTH AND MORE  
CONTAIN AT LEAST 50% OF THE  
ASSEMBLED BASES.

# Average contig length

---

## Average unigene length

---

### Contig N50 length

---

#### Unigene N50 length

= LENGTH AT WHICH  
THIS LENGTH  
CONTAIN AT LEAST  
ASSEMBLED

ique annotations in singlet

Unique annotations in singletons

---

Unique annotations in contigs

---

Unique annotations in unigenes

---

Average contig OHR

---

Average unigene OHR

ORTHOLOGOUS HIT RATE

# BASES IN ALIGNED

LENGTH OF BEST

(TRANSCRIPT COMPLETE)

Contig RBH count

# ORTHOLOGOUS HIT RATIO (OHR)

# BASES IN ALIGNED REGION

---

LENGTH OF BEST HIT

(TRANSCRIPT COMPLETENESS)

**Table 1 Metric trends and consistency**

| Metric                                                                                                     | Trend by         | Consistent over   | Trend by    | Consistent over | Fully consistent |
|------------------------------------------------------------------------------------------------------------|------------------|-------------------|-------------|-----------------|------------------|
|                                                                                                            | sequencing depth | sequencing depths | read length | read lengths    | metric           |
| Contig count                                                                                               | ↗                | ✓                 | ↘           | ✓               | ✗                |
| % of reads used in contigs<br>= <small>* BASES IN "GENES"<br/>+ TOTAL ASSEMBLED BASES</small>              | ↗                | ✓                 | ↗           | ✓               | ✓                |
| BP in contigs                                                                                              | ↗                | ✓                 | ↗           | ✓               | ✓                |
| % BP in contigs                                                                                            | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Average contig coverage                                                                                    | ↗                | ✗                 | ↗           | ✗               | ✗                |
| Average unigene coverage                                                                                   | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Contig read count COV                                                                                      | ↔                | ✓                 | ↘           | ✓               | ✗                |
| Unigene read count COV                                                                                     | ↔                | ✓                 | ↘           | ✓               | ✗                |
| Average contig length                                                                                      | ↗                | ✗                 | ↗           | ✗               | ✗                |
| Average unigene length                                                                                     | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Contig N50 length                                                                                          | ↗                | ✗                 | ↗           | ✗               | ✗                |
| Unigene N50 length<br>= <small>* NUMBER OF READS OF LENGTH EQUAL TO OR LONGER THAN THE N50 LENGTH</small>  | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Unique annotations in singletons                                                                           | ↔                | ✓                 | ↘           | ✓               | ✗                |
| Unique annotations in contigs                                                                              | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Unique annotations in unigenes                                                                             | ↗                | ✗                 | ↔           | ✓               | ✗                |
| Average contig OHR                                                                                         | ↗                | ✗                 | ↗           | ✗               | ✗                |
| Average unigene OHR<br>= <small>* NUMBER OF READS OF LENGTH EQUAL TO OR LONGER THAN THE N50 LENGTH</small> | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Contig RBH count                                                                                           | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Unigene RBH count                                                                                          | ↗                | ✓                 | ↗           | ✓               | ✓                |
| % of Annotated contigs with RBHs                                                                           | ↔                | ✗                 | ↗           | ✗               | ✗                |
| % of Annotated unigenes with RBHs                                                                          | ↔                | ✓                 | ↗           | ✗               | ✗                |
| Average contig CF                                                                                          | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Average unigene CF                                                                                         | ↗                | ✗                 | ↗           | ✗               | ✗                |
| Unique reverse annotations in contigs                                                                      | ↗                | ✓                 | ↗           | ✓               | ✓                |
| Unique reverse annotations in unigenes                                                                     | ↗                | ✓                 | ↗           | ✓               | ✓                |

**METHOD****Open Access**

# Evaluation of *de novo* transcriptome assemblies from RNA-Seq data

Bo Li<sup>1†</sup>, Nathanael Fillmore<sup>2†</sup>, Yongsheng Bai<sup>3</sup>, Mike Collins<sup>4</sup>, James A Thomson<sup>4,5,6</sup>, Ron Stewart<sup>4</sup> and Colin N Dewey<sup>2,7\*</sup>

<http://deweylab.biostat.wisc.edu/detonate/rsem-eval.html>

**"reference-free de novo transcriptome assembly evaluator"**

# RESULTS SO FAR...

| Assembly number | Trim quality (PHRED) | Normalized | Number of transcripts | Number of genes | unigene N50 |
|-----------------|----------------------|------------|-----------------------|-----------------|-------------|
| 1               | 20                   | Y          | 169 089               | 129 917         | 1655        |
| 2               | 10                   | Y          | 168 649               | 129 956         | 1640        |
| 3               | 5                    | Y          | 168 702               | 129 894         | 1664        |
| 4               | 2                    | Y          | 168 797               | 129 739         | 1660        |
| 5               | 20                   | N          | 174 531               | 144 326         | 1435        |
| 6               | 10                   | N          | 173 373               | 145 460         | 1425        |
| 7               | 5                    | N          | 175 339               | 145 410         | 1426        |
| 8               | 2                    | N          | 175 279               | 145 407         | 1427        |

# reads = 19 370 403

# assembled bases = 4 682 256 300