

A survey of popular ‘omics’ assembly tools

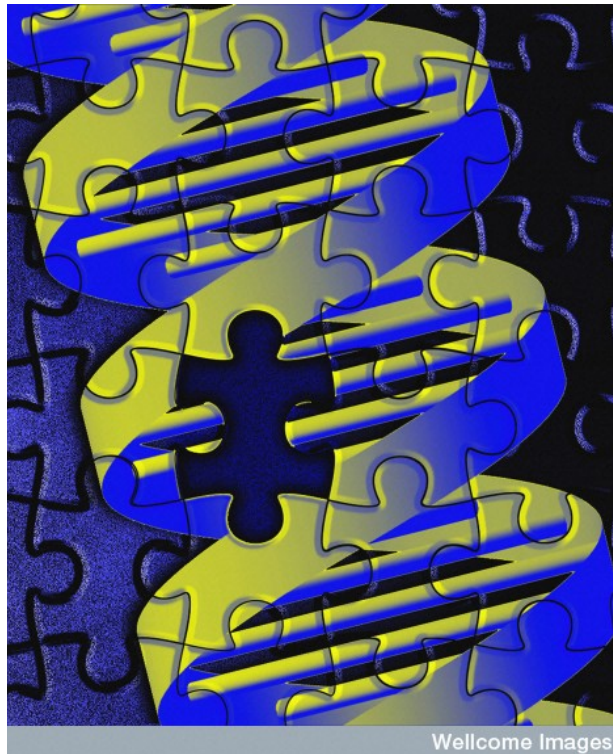


Image from [Wellcome Images](#)

Keith Bradnam
UC Davis Genome Center
July 2014

This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.



Overview

Between June 24th and July 1st 2014, I conducted a survey (using [an online Google Form](#)) to determine what tools are currently popular in the field of ‘omics’ assembly. The survey consisted of five questions, none of which were mandatory:

1. What is your preferred primary tool for eukaryotic genome assembly?
2. What is your preferred primary tool for bacterial/archaeal genome assembly?
3. What is your preferred primary tool for transcriptome assembly?
4. What is your preferred primary tool for metagenomics assembly?
5. How often do you generate genome/transcriptome/metagenome assemblies?

For the first four questions, I populated the survey form with a few examples of assembly tools but provided an option for ‘Other’. During the survey period, I would regularly update the form to include choices that people had provided in the ‘Other’ category. For the last question, five choices were available:

- I haven't yet made an assembly
- Maybe 1–2 a year
- Maybe 1–2 a month
- Maybe 1–2 a week
- More than 1–2 a week

Although it is common to use more than one piece of software as part of a sequence assembly pipeline, the goal of this survey was to identify the *primary* assembly tools that are *currently* being used. This survey was promoted via my [ACGT blog](#) and via [twitter](#).

Results

In the following tables, the names of assembly tools link to source web/FTP sites where the assembler software/code is available to download. For some commercial tools, the software is not directly available, but I try to link to the webpage that best describes the software.

Following the names of the software, I indicate a year that refers to the earliest date associated with the release of the software. This will sometimes correspond with the date of a publication that describes the software, but sometimes this date might precede the publication data if it is clear that the software was previously available (e.g. from developers website).

Other information is included below the tables as footnotes.

Table 1: Primary assembly tools for eukaryotic genome assembly

Assembler	Votes
<u>ALLPATHS-LG</u> (2010)	16
<u>ABYSS</u> (2008)	13
<u>SOAPdenovo</u> (2009)	10
<u>Velvet</u> (2007)	10
<u>SPAdes</u> (2012)	7
<u>Ray</u> (2010)	6
<u>Celera</u> (2004)	5
<u>CLC</u> (2008)	5
<u>MaSuRCA</u> (2012)	5
<u>Newbler</u> (2008) ¹	5
<u>Meraculous</u> (2011) ²	4
<u>SGA</u> (2010)	4
<u>Other</u> ³	2
<u>HGAP3</u> (2013)	2
<u>SeqMan NGen</u> (2007?) ⁴	1
<u>Minia</u> (2012)	1
<u>Platanus</u> (2012)	1
<u>Sprai</u> (2013)	1

¹ There is not really an official Newbler website, so readers are advised to check out Lex Nederbragt's [blog devoted to Newbler](#)

² Note, this is a link to an FTP site (ftp://ftp.jgi-psf.org/pub/JGI_data/meraculous/)

³ Two respondents indicated multiple tools (Celera/PBcR & ABYSS/SOAPdenovo/Platanus/Newbler)

⁴ This item was listed as 'DNASTAR' on the survey form (DNASTAR is the parent company).

Table 2: Primary assembly tools for bacterial/archaeal genome assembly

Assembler	Votes
<u>SPAdes</u> (2012)	29
<u>Velvet</u> (2007)	14
<u>ABYSS</u> (2008)	6
<u>ALLPATHS-LG</u> (2010)	6
<u>HGAP</u> (2013) ¹	6
<u>MIRA</u> (1998?) ²	6
<u>SOAPdenovo</u> (2009)	4
<u>CLC</u> (2008)	3
<u>Meraculous</u> (2011) ³	3
<u>SGA</u> (2010)	3
<u>IDBA-UD</u> (2011)	2
<u>A5</u> (2011)	1
<u>Celera</u> (2004)	1
<u>MaSuRCA</u> (2012)	1
<u>Newbler</u> (2008) ⁴	1
<u>Sprai</u> (2013)	1

¹ Five people responded with ‘HGAP2’ and one person reported ‘HGAP’

² Five people responded with ‘MIRA’ and one person reported ‘MIRA 4’

³ Note, this is a link to an FTP site (ftp://ftp.jgi-psf.org/pub/JGI_data/meraculous/)

⁴ There is not really an official Newbler website, so readers are advised to check out Lex Nederbragt’s [blog devoted to Newbler](#)

Table 3: Primary assembly tools for transcriptome genome assembly

Assembler	Votes
<u>Trinity</u> (2011)	57
<u>SOAPdenovo-Trans</u> (2011)	5
<u>Trans-ABYSS</u> (2010)	4
<u>CLC</u> (2008)	3
<u>MIRA</u>	2
<u>Oases</u> (2010)	2
<u>SeqMan NGen</u> (2007?) ¹	1
<u>IDBA</u> (2010) ²	1
IGC ³	1
Pertran ³	1
<u>Rnnotator</u> (2010)	1
Other ⁴	1

¹ This item was listed as ‘DNASTAR’ on the original survey (this is the parent company).

² These were survey responses that may have intended to be submitted as IDBA-Tran (2012), a transcriptome assembler.

³ I was unable to find a transcriptome assembler matching this name.

⁴ One person responded with ‘Velvet/Oases’

Table 4: Primary assembly tools for metagenome assembly

Assembler	Votes
<u>Velvet</u> (2007) ¹	17
<u>Ray</u> (2010) ²	11
<u>Ray-META</u> (2012) ²	6
<u>SeqMan NGen</u> (2007?) ³	2
<u>IDBA-UD</u> (2011) ⁴	2
<u>META-IDBA</u> (2011) ⁴	2
<u>Parallel-META</u> (2011)	2
<u>BAMBUS 2</u> (2011)	1
<u>Celera</u> (2004)	1
<u>CLC</u> (2008)	1
<u>Genovo</u> (2011)	1
<u>Meta is</u> ⁵	1
<u>MetaCortex</u> (2013)	1

¹ It is possible that some of these respondents intended to specify MetaVelvet (designed for metagenomics) rather than just Velvet.

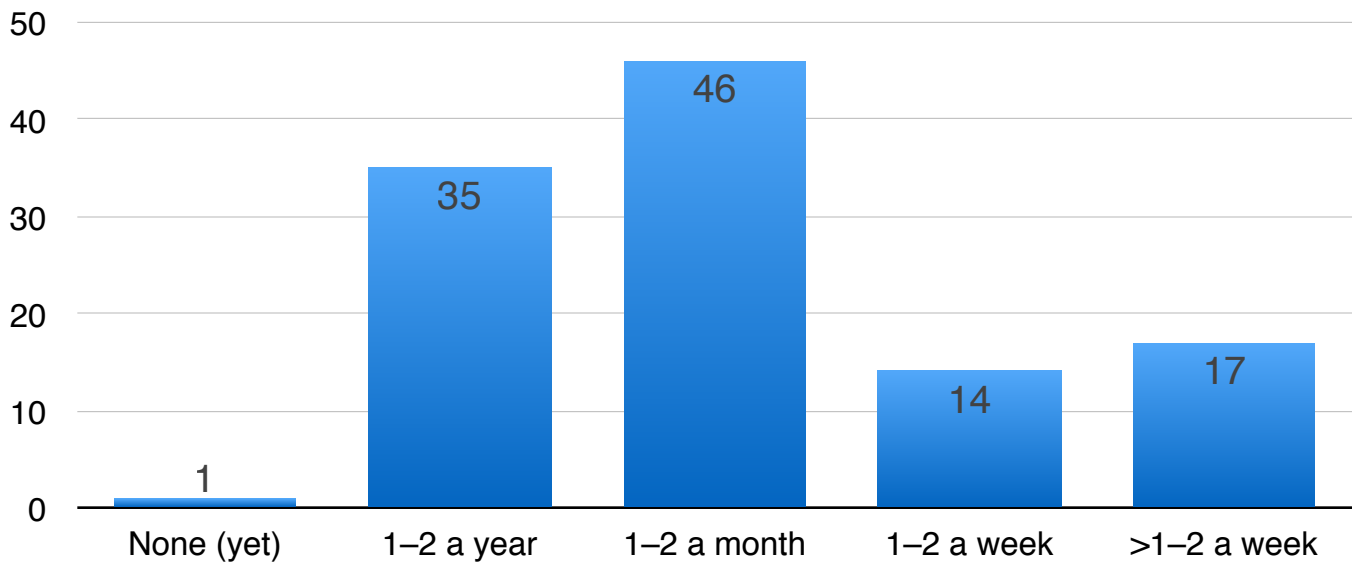
² It is possible that some respondents intended to specify Ray-META (designed for metagenomics) but instead chose Ray.

³ This item was listed as ‘DNASTAR’ on the survey form (DNASTAR is the parent company).

⁴ It is possible that some respondents intended to specify ‘META-IDBA’ (designed for metagenomics) but instead chose IDBA-UD.

⁵ I was unable to find a transcriptome assembler matching this name.

Figure 1: Frequency of assemblies generated by survey respondents



Discussion

The results of this informal survey show that there is a great deal of variety as to which ‘omics’ assemblers are currently in vogue. Of all these ‘omics’ disciplines, it is transcriptome assembly where one tool, Trinity, most clearly outranks all others in terms of popularity. Metagenomics is currently dominated by two tools (Velvet and Ray/Ray-Meta), as is the case with bacterial/archaeal assembly (SPAdes and Velvet).

It is only the discipline of eukaryotic genome assembly where several tools compete for preeminence. The four most popular assembly tools in this category are relatively old (the newest tool is from 2010) compared to tools in the other categories.

It is important to conclude this survey by noting that **popularity may not equate to quality**. The sequence assembly tools which are currently the most popular may simply be those that:

1. Have received more attention in the scientific literature
2. Work with sequencing data from more platforms than other tools
3. Have more applications than other tools (e.g. genome *and* metagenome assembly)
4. Are easier to install
5. Can be run on the available computational resources (some assembly tools require large amounts of memory and/or many CPUs)