

Small World Asynchronous Parallel Model for Genome Assembly

Jintao Meng, Jianrui Yuan, Jiefeng Cheng, Yanjie Wei, Shengzhong Feng

► **To cite this version:**

Jintao Meng, Jianrui Yuan, Jiefeng Cheng, Yanjie Wei, Shengzhong Feng. Small World Asynchronous Parallel Model for Genome Assembly. 9th International Conference on Network and Parallel Computing (NPC), Sep 2012, Gwangju, South Korea. pp.145-155, 10.1007/978-3-642-35606-3_17. hal-01551319

HAL Id: hal-01551319

<https://hal.inria.fr/hal-01551319>

Submitted on 30 Jun 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Small World Asynchronous Parallel Model for Genome Assembly

Jintao Meng^{1,2,4}, Jianrui Yuan^{2,3}, Jiefeng Cheng², Yanjie Wei^{2*}, Shengzhong Feng^{2*}

¹ Institute of Computing Technology, CAS, Beijing, 100190, PR.China

^{2*} Shenzhen Institutes of Advanced Technology, CAS, Shenzhen, 518055, PR. China

³ Central South University, Changsha, 410083, PR. China

⁴ Graduate University of Chinese Academy of Sciences, Beijing, 100049, China
{jt.meng, jr.yuan, jf.cheng, yj.wei, sz.feng}@siat.ac.cn

Abstract. Large de bruijn graph based algorithm is widely used in genome assembly and metagenetic assembly. The scale of this kind of graphs - in some cases billions of vertices and edges - poses challenges to genome assembly problem. In this paper, a one-step bi-directed graph is used to abstract the problem of genome assembly. After that small world asynchronous parallel model (SWAP) is proposed to handle the edge merging operation predefined in the graph. SWAP aims at making use of the locality of computing and communication to explore parallelism for graph algorithm. Based on the above graph abstraction and SWAP model, an assembler is developed, and experiment results shows that a factor of 20 times speedup is achieved when the number of processors scales from 10 to 640 when testing on processing *C.elegans* data.

Keywords: parallel computing, De Bruijn graph, genome assembly

1 Introduction

Current sequencing technology (Illumina Solexa [1], Applied Biosystems SoLiD[2], and Helicos Biosciences Heliscope[3]) allows one to read millions of short 35 to 100 nucleotide sequences per hour. Due to experimental errors, gaps, and genomic repeats, a much higher coverage depth of 50-fold to 300-fold is needed for accurate assembly. These factors contribute to a 300-fold to 1000-fold increase in the number of reads, which means there are billions of reads need to be processed, and this significantly complicate the genome assembly problem.

De Bruijn Assembler based on de-bruijn graph strategy [4,5] is well suitable for the current generation high throughput short reads assembly. In De Bruijn graph each vertex represents a length- k substring (k -mer) in a length- L read or its reverse complement. A directed edge connects two vertex u and v , if the $k-l$ length suffix of u is the same as the $k-l$ length prefix of v . Each input read is a path in the graph. By connecting such vertex pairs through edges, this approach will output the longest path without any branches as contigs. We denote the assemblers using De Bruijn graph

strategy as De Bruijn assembler.

The first De Bruijn assembler, EULER assembler [5] was proposed by Pevzner, who had transformed the fragment assembly problem into a variation of the classical Eulerian path problem by dividing reads into k -mers and then constructing k -mers into a path graph. This opens new possibilities for repeat resolution and generating error-free solutions of the large-scale fragment assembly problem. Programs such as Velvet[6], SOAPdenovo[7], and IDBA[8] implicitly use this framework but are slightly different in local details. Velvet manipulates these De Bruijn graphs efficiently to both eliminate errors and resolve repeats by error correction algorithm. SOAPdenovo implement pre-assembler error correction on human genome assembly, after this operation the proportion of error free reads was improved from 64% to 70%, and nearly 60% percent of k -mers was filtered from the graph. IDBA also adopt pre-assembler error filtering technique, which can save nearly 40-80% of memory compared with velvet. The second feature of IDBA is that it iterates from small k -mer to large k -mer to get longer contigs. So the quality of contigs is better than other tools.

The above assemble tools can only run on single machine, the human genome assembly with current sequencing technology needs about 1TB memory and takes weeks or even months on single server. The situation will be even worse for larger genome assembly or meta-genome assembly.

Parallel algorithm for sequencing assembly is an alternative to solve the problem. Parallel assemblers included ABySS[9] and YAGA[10-12], are both based on De Bruijn graph strategy. ABySS distributes k -mers to multi-servers to build a distributed De Bruijn graph, and error removal and vertex merging were implemented over MPI communication messages. YAGA constructs the distributed bi-directed De Bruijn graph by maintaining edge tuples in a community of servers. Unanimous chain compaction problem in YAGA was transformed to undirected list ranking, and then the authors designed a modified sparse ruling set algorithm for undirected lists. The computational complexity of YAGA is given by $O(n/p)$ compute time, $O(n/p)$ communication volume, and $O(\log^2 n)$ communication rounds, where n is the number of nucleotides in all reads, p denotes the number of processors.

Efficient and scalable frameworks or libraries for distributed graphs are essential to parallel assembly based on De Bruijn graph. Existing works, such as BSPLib [13-15], CGMgraph [16], PBGL [17,18], Prejel [19], are based on BSP [20] model. The BSP model has advantage on simple computation-communication programming model, whereas the barrel principle exists in the computation-communication phase and the synchronous phase over large clusters limits the scalability of the model. To our knowledge, the scalability these implementations under BSP model has not been evaluated beyond several hundreds of computers [19]. No genome assembly tools have adapted the BSP library, although YAGA has used the BSP idea in its design on parallel list ranking algorithm implicitly. Another parallel programming model, MapReduce [21], has strength in loosely coupled work such as frequency statistics, sorting, indexing, and machine learning etc, and these works can be easily distributed to clusters. However graph algorithm is a tight coupled work and dividing one graph into several meaningful sub-graphs is still a challenging problem.

This paper first demonstrates a one-step bi-directed graph for the problem of genome assembly. Genome can be recovered by merging semi-extended edges to full-extended edges. Then small world asynchronous parallel (SWAP) model is proposed

to realize edge merging over a distributed one-step bi-directed graph. Specially, we implement an assembler using the SWAP model. Given the number of processes p , the complexity of this problem is reduced to $O(n/p)$ parallel compute time, $O(n/p)$ communication round, and $O(n \log(n)/p)$ communication volume, here n is total length of input sequences. Simulation shows that Assembler has a factor of 20 times speedup when the number of processors scales from 10 to 640.

The rest of the paper is organized as follows: Section 2 abstracts the De Bruijn graph based genome assembly problem; Section 3 describes the SWAP model for large scale graphs with small world property, then an assembler, as SWAP's first application, is illustrated. Experimental results will be present in section 4. Finally section 5 concludes this paper.

2 Abstraction of De Bruijn Assembly

Let $s \in M^l$ be a string of length L , where $M = \{a, t, c, g\}$. Any substring $\alpha = s[j]s[j+1] \dots s[j+k-1]$, $0 \leq j < L-k+1$ is a k -mer of s . The set of all k -mers of a given string s is written as $\mathbb{Z}(s, k)$, here k must be odd. The reverse complement of a k -mer α , denoted by α' , is obtained by reversing α and complementing each base ($\alpha'[i] = \alpha[k-i+1]'$) by the following bijection of $\Sigma, \bar{\Sigma}: \{a \rightarrow t, t \rightarrow a, c \rightarrow g, g \rightarrow c\}$. Note that $\alpha[i] = \alpha[i]''$ and $\alpha = \alpha''$.

A k -molecule $\hat{\alpha}$ is a pair of complementary k -mers $\{\alpha, \alpha'\}$. Let \geq be the partial ordering relation among the string of equal length such that $\alpha \geq \beta$ indicates that the string α is lexicographically larger than β . We designate the lexicographically larger of the two complementary k -mers as the positive k -mer, denoted as α^+ , and the lexicographically smaller one as the negative k -mer, denoted as α^- , here $\alpha^+ \geq \alpha^-$. We choose the positive k -mer α^+ as the representative k -mer $\hat{\alpha}$ of the k -molecule. The set of all k -molecules of a given string s is called the k -spectrum of s and is written as $\mathbb{S}(s, k)$. Noted that $\mathbb{S}(s, k) = \mathbb{S}(s', k)$.

The notation $su f(a, l)$ ($pre(a, l)$, respectively) is used to denote the length l suffix (prefix, respectively) of string a . Let the symbol \circ denotes the concatenation operation between two strings, and the number of edges attached to k -molecule $\hat{\alpha}$ is denoted as $degree(\hat{\alpha})$. The number of edges pointing out from k -molecule $\hat{\alpha}$ is denoted as $degree(\hat{\alpha})$.

Definition 1. The vertex set V_s is defined as k -spectrum of s ,

$$V_s = \mathbb{S}(s, k) \quad (1)$$

Definition 2. The 1-step bi-directed edge set E_s^1 is defined as follows:

$$E_s^1 = \{e_{\alpha\beta}^1 = (\alpha, \beta, d_\alpha, d_\beta, c_{\alpha\beta}^1) \mid \forall \hat{\alpha}, \hat{\beta} \in \mathbb{S}(s, k), \text{ suf}(\alpha, k-1) = \text{pre}(\beta, k-1) \wedge (\alpha \circ \beta[k-1] \in (\mathbb{Z}(s, k+1) \vee \mathbb{Z}(s', k+1)))\} \quad (2)$$

Equations (2) declares that any two overlapped k -molecules can be connected with a 1-step bi-directed edge, if they are continuous in sequence s or its complementary.

Here d_α is the direction of k -mer α , if $\alpha = \alpha^+$, $d_\alpha = '+'$, otherwise $d_\alpha = '-'$.

Set $c_{\alpha\beta}^1$ is initialized with one element $\beta[k-1]$, and $\text{suf}(\alpha \circ c_{\alpha\beta}^1, k) = \beta$.

Property 1. Given two k -molecules $\hat{\alpha}, \hat{\beta} \in \mathbb{S}(s, k)$, there will be four possible connections, and for each type of connection exactly two equivalent 1-step bi-directed edge exist,

1. $e_{\alpha^+\beta^+}^1 = (\alpha^+, \beta^+, +, +, c_{\alpha^+\beta^+}^1), e_{\alpha^-\beta^-}^1 = (\alpha^-, \beta^-, -, -, c_{\alpha^-\beta^-}^1),$
2. $e_{\alpha^+\beta^-}^1 = (\alpha^+, \beta^-, +, -, c_{\alpha^+\beta^-}^1), e_{\alpha^-\beta^+}^1 = (\alpha^-, \beta^+, -, +, c_{\alpha^-\beta^+}^1),$
3. $e_{\alpha^-\beta^+}^1 = (\alpha^-, \beta^+, -, +, c_{\alpha^-\beta^+}^1), e_{\alpha^+\beta^-}^1 = (\alpha^+, \beta^-, +, -, c_{\alpha^+\beta^-}^1),$
4. $e_{\alpha^-\beta^-}^1 = (\alpha^-, \beta^-, -, -, c_{\alpha^-\beta^-}^1), e_{\alpha^+\beta^+}^1 = (\alpha^+, \beta^+, +, +, c_{\alpha^+\beta^+}^1).$

In each type of connection, the first bi-directed edge and the second one correspond to the same bi-directed edge, but in different form. Within a distributed edge representation situation, the first bi-directed edge in each type will be attached with k -molecule $\hat{\alpha}$, and the second one will be with $\hat{\beta}$. Figure (1) illustrates four possible connections and examples of a 1-step bi-directed edge graph.

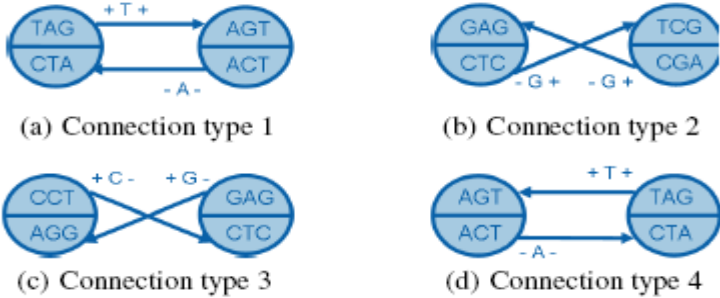


Fig. 1. The illustration of four possible connections.

Definition 3. 1-step bi-directed de bruijn graph of order k for a string s can be:

$$G_k^1(s) = \{V_s, E_s^1\} \quad (3)$$

Definition 4. Given two 1-step bi-directed edge $e_{\alpha\beta}^1 = (\alpha, \beta, d_\alpha, d_\beta, c_{\alpha\beta}^1)$ and $e_{\beta\gamma}^1 = (\beta, \gamma, d_\beta, d_\gamma, c_{\beta\gamma}^1)$, if $e_{\alpha\beta}^1.d_\beta = e_{\beta\gamma}^1.d_\beta$ and $\text{degree}(\hat{\beta})=2$, we can get 2-step bi-directed edge $e_{\alpha\gamma}^2 = (\alpha, \gamma, d_\alpha, d_\gamma, c_{\alpha\gamma}^2)$ by merging $e_{\alpha\beta}^1$ and $e_{\beta\gamma}^1$. Here $c_{\alpha\gamma}^2 = c_{\alpha\beta}^1 \circ c_{\beta\gamma}^1$. Let the symbol \oplus denote **edge merging operation** between two

bi-directed edges attached to one k -molecule, then **edge merging operation** can be written as,

$$e_{\alpha\gamma}^2 = e_{\alpha\beta}^1 \oplus e_{\beta\gamma}^1 \quad (4)$$

$$\text{or } e_{\gamma\alpha}^2 = e_{\gamma\beta}^1 \oplus e_{\beta\alpha}^1 \quad (5)$$

According to property 1, equation (4) and equation (5) correspond to one edge merging operation. Then z -step bi-directed edge can be defined as:

$$e_{\alpha\gamma}^z = e_{\alpha\beta}^x \oplus e_{\beta\gamma}^y, \text{ iff } \exists \beta, e_{\alpha\beta}^x \cdot d_\beta = e_{\beta\gamma}^y \cdot d_\beta, \text{ deg } \text{ree}(\widehat{\beta}) = 2, z = x + y \quad (6)$$

Definition 5. Given an n -step bi-directed edge $e_{\alpha\beta}^m = (\alpha, \beta, d_\alpha, d_\beta, C_{\alpha\beta}^m)$, if k -molecule α or β has only one another bi-directed edge $e_{\gamma\alpha}^t = (\gamma, \alpha, d_\gamma, d_\alpha, C_{\gamma\alpha}^t)$ or $e_{\beta\gamma}^t = (\beta, \gamma, d_\beta, d_\gamma, C_{\beta\gamma}^t)$ respectively, then $e_{\alpha\beta}^m$ can be extended by $e_{\gamma\alpha}^t$ or $e_{\beta\gamma}^t$, we regard this edge as semi-extended edge, and the corresponding k -molecule as semi-extended vertex. If $e_{\alpha\beta}^m$ cannot be extended by any edge, we call this edge a full-extended edge.

Given a set of string or reads $S = \{s_1, s_2, \dots, s_h\}$, a one-step bi-directed De Bruijn graph of S with order of k is $G_k^1(S) = \{V_S, E_S^1\} = \left\{ \bigcup_{1 \leq i \leq h} V_{s_i}, \bigcup_{1 \leq i \leq h} E_{s_i}^1 \right\}$. The key property of this bi-directed De Bruijn graph $G_k^1(S)$ is that each read can be recovered by traversing the corresponding path in either direction, concatenating $(k-1)$ -molecule prefix of the first node and the edge labels on the path. As all input reads of assembler are derived from chromosomes, each chromosome can now be seen as a long path in this graph. However because of read errors, and repeats in the sequence, we cannot expect to see continuity in sampling, our goal is to recover the genome as a large set of contigs by merging semi-extended edges into full-extended edges.

3 Assembler over SWAP

Vertices in large scale real world graph (such as social network, web link graph, et) always have limited number of neighbors, little computing work, and constant number of edges randomly connected to other vertices, this phenomenon is denoted as small world property. For a given vertex, its small world includes all its edges, neighbors and itself. Then any computing and communication work of a vertex can be done in its small world. As long as the work of a vertex on a graph does not interrupt others, we can run computational work of those vertices in parallel. Here we will present our work on pursuing parallelism in the computation of bi-directed graph for genome assembly.

Inspired from CSMA/CA in wireless networks [22], Small World Asynchronous Parallel model (SWAP) aims to improve parallelism on processing large scale graph problem with small world property. After having distributed graph over a network of

processors, the main schedule of SWAP can be defined as a combination of following three steps:

1. **Lock** operation is applied to each vertex's small world, which includes itself and its neighbors.
2. **Computation** and modification will be performed in each vertex's small world..
3. **Unlock** operation will be triggered after each computation step.

The basic schedule of SWAP is Lock-Computation-Unlock. Because of the locality of computing and communication in the small world, SWAP model utilizes local synchronization and global asynchronization mechanism to maximize underline parallelism for the graph algorithm.

An assembler over SWAP is the first application using SWAP model. In the following paragraphs we will describe its data structure on distributed de bruijn graph, strategy on error removal, and the edge merging algorithm, respectively.

3.1 Parallel constuction of distributed one-step bi-directed De Bruijn graph

$S = \{s_1, s_2, \dots, s_m\}$ is m sequences sampled from a genome of length g , the total length of all these sequence is n . we aim to construct a one-step bi-directed De Bruijn graph $G_k^1(S)$ with $O(n)$ vertexes and edges distributed among p processors such that each processor store $O(n/p)$ vertices.

Input sequences can be broken into overlapping k -molecules by sliding a window of length k along the input sequence. Each processor maintains a hash table to store k -molecules, and each k -molecule is represented as a base-4 number of its positive k -mer. Numerical values $\{0,1,2,3\}$ are assigned to bases $\{A,C,G,T\}$. The location of a given k -molecule can be computed by take the mod of a large prime number and then take the mod of the number of processors. The large prime number is used to evenly distribute k -molecules to all processors.

A single k -molecule can have up to eight edges, and each of them corresponds to a one-base extension, $\{A, C, G, T\}$ in either direction. The adjacent k -molecule can be easily generated by adding the base extension in the edge set to the source k -molecule.

The construction of one-step bidirected De Bruijn graph can be achieved in $O(n/p)$ parallel compute time, $O(1)$ round of all-to-all communication, and $O(n/p)$ parallel communication volume.

3.2 Error removal

Sequencing errors make the assembly problem more complex. To identify errors, we assume that the errors are random, and they are unlikely to occur twice in the same base. As each base in the genome is sampled on average as many times as the overage

number, the erroneous k -molecule will have lower frequency compared to the correct ones. According to this principle, we identify all k -molecules with low frequency as erroneous k -molecules, and delete all of them from our vertex set of the graph. The complexity of this step is $O(n/p)$ parallel compute time.

3.3 Edge merging

One step bidirected graph generated in the previous section will likely have many long chains, and each corresponds to a sequence that can be unambiguously assembled into a single contig. We will merge these chains into full-extend edges using **Algorithm 1**.

Algorithm 1 Edge Merging Operation Algorithm

Iteration:

1: Element selection operation

For each semi-extended k -molecule $\hat{\alpha}$ in V_i , if $\hat{\alpha}$'s two neighbors $\hat{\beta}$ and $\hat{\gamma}$ are connected by edge $e_{\beta\alpha}^u$ and $e_{\alpha\gamma}^v$, then $e_{\beta\alpha}^u$ and $e_{\alpha\gamma}^v$ can be merged as $e_{\beta\gamma}^{u+v} = e_{\beta\alpha}^u \oplus e_{\alpha\gamma}^v$, $e_{\gamma\beta'}^{u+v} = e_{\gamma\alpha'}^v \oplus e_{\alpha'\beta'}^u$

2: Lock

In order to merging $e_{\beta\alpha}^u$, $e_{\alpha\gamma}^v$, we need to send Lock messages to lock k -molecules $\hat{\alpha}$, and its neighbours $\hat{\beta}$, $\hat{\gamma}$.

3: Computing

$e_{\beta\alpha}^u$, $e_{\alpha\gamma}^v$ and $e_{\gamma\alpha'}^v$, $e_{\alpha'\beta'}^u$ will be merged into one edge $e_{\beta\gamma}^{u+v}$. That means the original two edges $e_{\beta\alpha}^u$, $e_{\alpha\gamma}^v$ will be deleted, and new edge $e_{\beta\gamma}^{u+v}$, $e_{\gamma\beta'}^{u+v}$ between $\hat{\beta}$ and $\hat{\gamma}$ will be added.

4: Unlock

k -molecule $\hat{\alpha}$ sends unlock messages to unlock k -molecule $\hat{\alpha}$, $\hat{\beta}$ and $\hat{\gamma}$.

Output:

Return Full-extended edges.

A subset of vertices and their associated edges are stored in processor i , and the set of semi-extended vertices stored in this processor is denoted as V_i . We aim to delete all semi-extended vertices and merge their associated semi-extended edges to form full-extended edges.

As the bi-directed graph $G_k^1(S)$ is distributed over p processors, each processor will store a subset of semi-extended k -molecules V_i and the average number of k -molecules in V_i is $O(n/p)$. Then the expected computational complexity of each processor on edge merging is given by $O(n/p)$ parallel compute time, $O(n/p)$ communication round, and $O(n \log(n)/p)$ communication volume.

4 Experimental Results

The assembler is written in C++ and MPI. The hardware and software architecture supporting this assembler is demonstrated in Fig 2. We use Dawning 5000 as high performance cluster, which has 40 16-core servers with 32GB memory. The distributed file system is lustre. All the components are interconnected with infiniband 20Gbit Router.

Perl scripts [23] are used to generate the following two theoretical datasets: 50x coverage of Yeast chromosomes containing 17 million reads, and 50x coverage of

C.elegans chromosomes containing 141 million reads. The error rate is set to be 1%, and the length of reads ranges from 36bp to 50bp. The primary goal of this experiment is to demonstrate the scalability of SWAP model on handling large-scale graphs using parallel system with distributed memory.

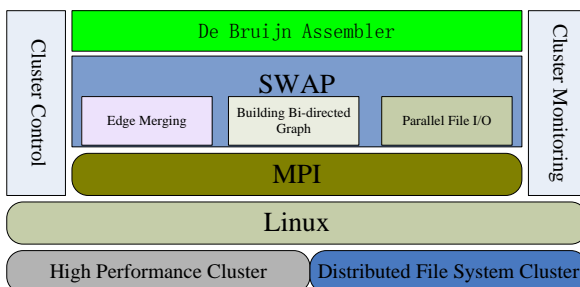


Figure 1. Hardware and software architecture of the assembler on SWAP.

We first test the performance of this assembler on Yeast dataset. The runtime of assembler is displayed in figure 3, and the time is divided into three phases, Parallel File I/O, graph building, and edge merging. The first phase is the time spent on reading dataset from a distributed file system, the second phase is the time used to construct the one-step bi-directed graph over the cluster, and the last phase is the time cost on edge merging operations. The run time is dominated by the third phase, where hundreds of processors are sending messages to lock their neighbors, merging edges, deleting semi-extended nodes and edges, and unlocking their neighbors. Figure 3 shows that this phase has good scalability. The speedup is about 50 when the number of processor scales from 10 to 640 and the overall runtime of assembler on Yeast dataset is reduced by a factor of 30.

The C.elegans dataset is nearly ten times larger than Yeast, and its corresponding data on time usage is demonstrated in figure 4. In this figure, the time used in parallel file I/O is very short compared to the other two phases. The graph construction phase has a decreasing trend on its running time. This phase have a speedup of 35x. This speedup is slightly smaller than that of the yeast dataset. The total speedup of all three phases on C.elegans dataset is about 20.

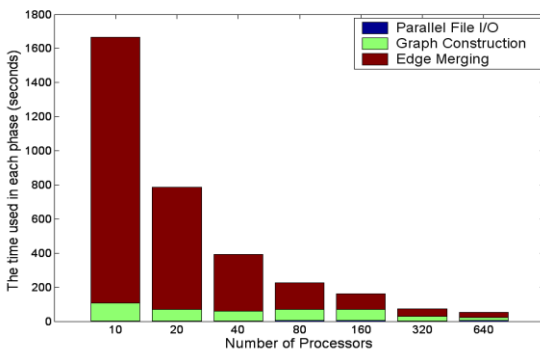


Fig.2 Time usage analysis in three phase on Yeast dataset

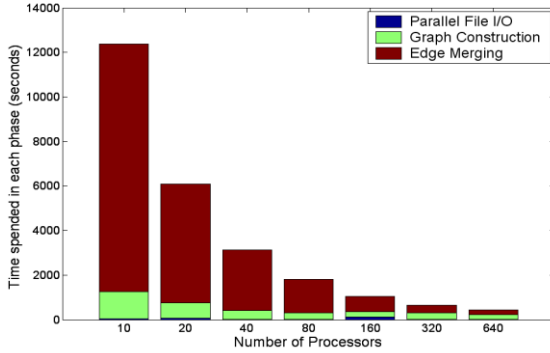


Fig.3 Time usage analysis in three phase on C.elegans dataset

4 Conclusion

In this paper, we abstracted the problem of genome assembly using De Bruijn strategy. By constructing the one-step bi-directed graph over k-spectrum of input sequences, the unanimous path compaction problem in generic genome assembly was transformed to merge semi-extended edges in our bi-directed graph, and the final contigs are full-extended edges in our method. The proposed SWAP introduced local synchronization and global asynchronization mechanism to maximize the parallelism in the graph algorithm. SWAP model applies the Lock-Computation-Unlock scheme to each vertex's small world. Based on SWAP model, we developed a De Bruijn assembler, and simulation results show that when the number of processors scales from 10 to 640, a factor of 30 and 20 speedup, can be achieved for assembling Yeast and C.elegans genomes, respectively.

Acknowledgements

This work is supported by NSFC (Grant No. 61103049) and Shenzhen Research Fund (Grant No.JC201005270342A). The author also thanks Bingqiang Wang from BGI, and Prof. Francis Y.L. Chin from HKU for their suggestions on this work.

References

1. S. Bennet. Solexa ltd, Pharmacogenomics, Vol.5, No.4, pp. 433-438, June. 2004.
2. V. Pandey, R.C. Nutter, E. Prediger. "Applied Biosystems SOLiDTM System: Ligation-Based Sequencing." Next Generation Genome Sequencing: Towards Personalized Medicine, Wiley, 2008.

3. Business Wire, "Helicos biosciences enters molecular diagnostics collaboration with renowned research center to sequence cancer-associated genes," Genetic Engineering and Biotechnology News, 2008.
4. R.M. Idury, M.S. Waterman, "A New Algorithm for DNA Sequence Assembly," *Journal of Computational Biology*, vol. 2, no. 2, pp. 291-306, 1995
5. P.A. Pevzner, H. Tang, M.S. Waterman, "An Eulerian path approach to DNA fragment assembly," *Proceedings of the National Academy of Sciences of the United States of America (PNAS)*, Aug. 2001, vol. 98, no. 17, pp. 9748-9753.
6. D.R. Zerbino, E. Birney, "Velvet: algorithms for de novo short read assembly using De Bruijn graphs," *Genome Research*, vol. 18, no.5, pp.821-829, 2008.
7. R. Li, H. Zhu, J. Ruan, W. Qian, X. Fang, Z. Shi, Y. Li, S. Li, G. Shan, K. Kristiansen, S. Li, H. Yang, J. Wang, J. Wang, "De novo assembly of human genomes with massively parallel short read sequencing" *Genome Research*, vol. 20, no. 2, pp. 265-272, 2010.
8. Y. Peng, Henry C.M. Leung, S.M. Yiu, Francis Y. L. Chin, "IDBA-A Practical Iterative De Bruijn Graph De Novo Assembler," in *RECOMB 2010*, vol. 6044, pp. 426-440.
9. J.T. Simpson, K. Wong, S.D. Jackman, J.E. et al, "ABySS: a parallel assembler for short read sequence data," *Genome Research*, vol. 19, no. 6, pp. 1117-1123, 2009
10. B.G. Jackson, S. Aluru, "Parallel Construction of Bidirected String Graphs for Genome Assembly," in *Proc. of the 37th International Conference on Parallel Processing (ICPP'08)*, Sept. 2008, pp. 346-353.
11. B.G. Jackson, P.S. Schnable, S. Aluru, "Parallel short sequence assembly of transcriptomes," *BMC Bioinformatics* 10(S-1), 2009.
12. B.G. Jackson, M. Regennitter, X. Yang, P.S. Schnable, S. Aluru, "Parallel de novo assembly of large genomes from high-throughput short reads," in *Proc. of the 24th International Symposium on Parallel & Distributed Processing (IPDPS'10)*, Atlanta, 2010
13. Richard Miller, "A Library for Bulk-Synchronous Parallel Programming," in *Proc. British Computer Society Parallel Processing Specialist Group Workshop on General Purpose Parallel Computing*, 1993.
14. M.W. Goudreau, K. Lang, S.B. Rao, T. Suel, and T. Tsantilas, "Portable and Efficient Parallel Computing Using the BSP Model," *IEEE Transactions on Computers*, vol. 48, no. 7, pp.670-689, 1999.
15. O. Bonorden, Ben H.H. Juurlink, Ingo von Otte, and I. Rieping, "The Paderborn University BSP (PUB) Library," *Parallel Computing*, vol. 29, no. 2, pp. 187-207, 2003.
16. A. Chan, F. Dehne, "CGMGRAPH/CGMLIB: Implementing and Testing CGM Graph Algorithms on PC Clusters and Shared Memory Machines," *International Journal of High Performance Computing Applications*, vol. 19, no. 1, pp. 81-97, 2005
17. D. Gregor, A. Lumsdaine, "The Parallel BGL: A Generic Library for Distributed Graph Computations," in *Proc. of Parallel Object-Oriented Scientific Computing (POOSC)*, 2005.
18. D. Gregor, A. Lumsdaine, "Lifting Sequential Graph Algorithms for Distributed-Memory Parallel Computation," in *Proc. of the 20th annual ACM SIGPLAN conference on object-oriented programming systems, languages, and applications (OOPSLA'05)*, pp. 423-437.
19. G. Malewicz, M.H. Austern, Aart J. C. Bik, J.C. Dehnert, I. Horn, N. Leiser, G. Czajkowski, "Pregel: a system for large-scale graph processing," In *SIGMOD'10 Proceedings of the 2010 international conference on Management of data*, 2010, pp. 135-146, New York.
20. Leslie G. Valiant, A bridging model for parallel computation, *Communications of the ACM*, vol. 33, no. 8, Aug. 1990.
21. J. Dean, S. Ghemawat, "MapReduce: simplified data processing on large clusters," *Communications of the ACM - 50th anniversary issue: 1958 - 2008* vol. 51, no. 1, 2008
22. Andrew S. Tanenbaum, *Computer Networks*, Prentice Hall, New Jersey, 2003
23. W. Zhang, J. Chen, Y. Yang, Y. Tang, J. Shang, B. Shen, "A practical comparison of de novo genome assembly software tools for next-generation sequencing technologies," *PLoS ONE*, vol. 6, no. 3. March 2011