

DATABASE

Open Access

ELM: enhanced lowest common ancestor based method for detecting a pathogenic virus from a large sequence dataset

Keisuke Ueno¹, Akihiro Ishii² and Kimihito Ito^{1*}

Abstract

Background: Emerging viral diseases, most of which are caused by the transmission of viruses from animals to humans, pose a threat to public health. Discovering pathogenic viruses through surveillance is the key to preparedness for this potential threat. Next generation sequencing (NGS) helps us to identify viruses without the design of a specific PCR primer. The major task in NGS data analysis is taxonomic identification for vast numbers of sequences. However, taxonomic identification via a BLAST search against all the known sequences is a computational bottleneck.

Description: Here we propose an enhanced lowest-common-ancestor based method (ELM) to effectively identify viruses from massive sequence data. To reduce the computational cost, ELM uses a customized database composed only of viral sequences for the BLAST search. At the same time, ELM adopts a novel criterion to suppress the rise in false positive assignments caused by the small database. As a result, identification by ELM is more than 1,000 times faster than the conventional methods without loss of accuracy.

Conclusions: We anticipate that ELM will contribute to direct diagnosis of viral infections. The web server and the customized viral database are freely available at <http://bioinformatics.czc.hokudai.ac.jp/ELM/>.

Keywords: Next generation sequencing, Virus discovery, Diagnostic virology, Virome, Taxonomic identification

Background

Most emerging infectious diseases are zoonoses, the pathogens of which are transmitted between humans and animals. The 2009 pandemic H1N1 influenza virus spread worldwide through reassortment that exchanged a gene segment between pigs and humans [1]. Recently, cases of influenza A virus H7N9 transmitted from birds to humans have been reported [2]. The 2003 severe acute respiratory syndrome (SARS) outbreak originated from the transmission of a novel bat coronavirus [3]. For the sporadically endemic Ebola virus, bats are suspected to be the natural reservoir, but this is still controversial [4]. Vector-borne zoonoses caused by transmission of viruses through mosquitoes and ticks have also become a public health concern. The 1999 outbreak of West Nile virus (WNV) that occurred in New York was caused by

the transmission of the WNV among birds, horses and humans via mosquitoes [5]. Similarly, severe fever with thrombocytopenia syndrome (SFTS) was found to be due to a virus transmitted by ticks [6].

To prepare for the risk of emerging infectious diseases, we need to identify pathogenic viruses through surveillance of livestock and wild animals. Although universal PCR primers against 16S ribosomal RNA are available for the identification of bacteria, we needed specific PCR primers to identify viruses. In recent years, NGS technologies have become available for identifying novel viruses that cannot be found by Sanger sequencing due to the difficulty of isolation and passage culture [7].

The taxonomic classification of metagenomic sequences is an important task in NGS data analyses [8]. It has been widely applied to investigate the relationship between human health and the microbiome [9]. Recently, a metagenomic analysis of the virome in a monkey infected with simian immunodeficiency virus was conducted, suggesting that the virome was associated with enteropathy caused

* Correspondence: itok@czc.hokudai.ac.jp

¹Division of Bioinformatics, Research Center for Zoonosis Control, Hokkaido University, Sapporo, Hokkaido 001-0020, Japan

Full list of author information is available at the end of the article

by HIV [10]. Through the first screening with NGS, the novel influenza virus H17N10 was identified in bats from metagenomic samples [11].

The taxonomic classification of NGS data uses sequence similarity searches such as BLASTX and BLASTN [12] to assign each sequence into a specific taxon based on the hits. However, with the similarity-based approach it is difficult to decide the resolution of assignments because the resolution depends on whether the sequences are conserved or species specific. The metagenome analyzer (MEGAN) employs the lowest common ancestor (LCA) concept in graph theory to estimate the taxonomical contents of samples [8]. MEGAN evaluates the resolution of similarity-based assignments as the level of taxonomy based on the LCA.

The LCA is the closest taxon shared among two or more taxa found by a BLAST search for a read. When multiple taxa are found by the BLAST search with sufficiently reliable BLAST scores, the common ancestor is a high-level taxon. The LCA assignments to high-level taxa are associated with conserved sequences. When a single taxon is found by a BLAST search for a read, the common ancestor still remains a low-level taxon. The LCA assignments to low-level taxa are associated with species-specific sequences. Thus, the LCA assignments to low-level taxa are more suitable for resolving closely related organisms than those to high-level taxa.

The SORT-ITEMS [13] and CARMA3 [14] methods extended the LCA using a reciprocal BLAST search to reduce false positives in assignments. CARMA3 introduced the concept of the mutation rate into the LCA algorithm, and reinforced the reciprocal BLAST search to identify a novel taxon, relatives of which are numbered [14].

While taxonomic classification of metagenomic sequences has been developed with respect to accuracy, NGS technologies continue to improve sequencing throughput, and require considerable computational time and resources to perform taxonomic classification. The throughput of Roche 454 sequencing is 700 Mb with an average length of 400–800 bases. The present throughputs of NGS have become over 1 Gb with Illumina sequences of 600 Gb and an average length of ~100 bases, SOLiD sequences of 20 Gb with an average length of ~50 bases, and Ion Torrent PGM sequences of 1 Gb with an average length of ~200 bases [7]. These massive sequencing data prevent the fast detection of infecting viruses from metagenomic samples.

To reduce the computational time, we constructed a customized database composed only of viruses for the BLAST search. However, customized databases also increase accidental hits, i.e. the match of host sequences to viral genomic sequences. Here, we introduce ELM with a customized viral database for taxonomic identification. The method is based on the assumption that

valid hits, the match of viral sequences to viral genomic sequences, raise the probability of finding other similar genomic sequences in the BLAST search. In other words, true assignments with the LCA should be sensitive to the threshold of the bit score in the BLAST search. Consequently, ELM can suppress the rise of false positive assignments while saving computational time and resources.

Construction and content

The ELM server performs taxonomic identification of viral sequences from NGS datasets via three steps (Figure 1). In step one, the server carries out a BLASTN search for a customized database of viral genomic sequences. In step two, the server performs the LCA-based taxonomic assignments using MEGAN software [8] with default parameters. In step three, the server iterates the LCA assignments with different parameters for the threshold of the bit scores for the BLAST hits and investigates the taxa in which the number of assigned reads is significantly changed. In this step, the server provides a novel criterion for evaluating the LCA assignments.

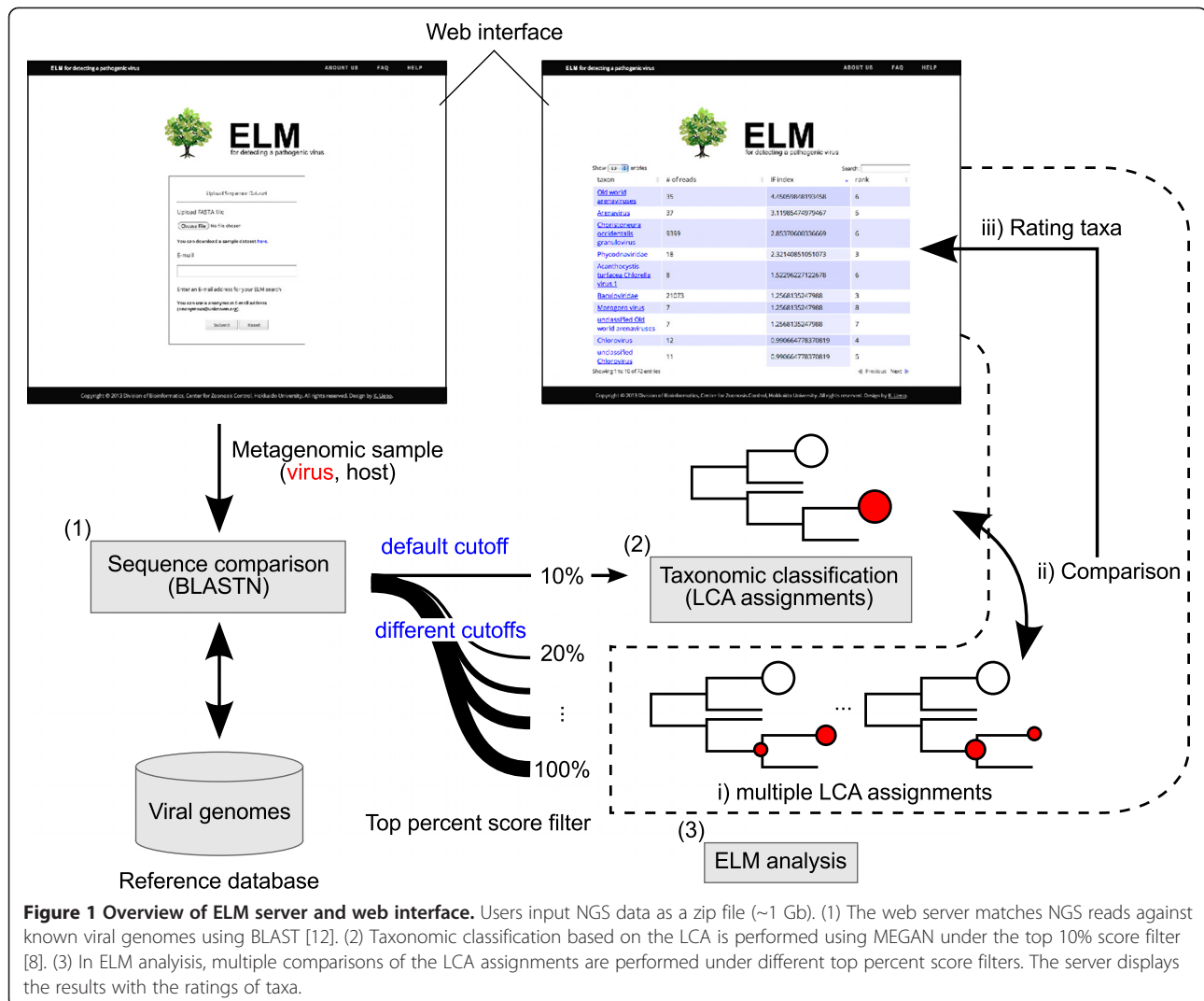
BLAST search for customized database

To reduce the computational time and save disk space, we constructed a customized database composed only of viral genomic sequences for a BLASTN search. First, the RefSeq genomic sequences were downloaded from the NCBI. Then a total of 3,336 viral genomic sequences were selected using a custom-made script program and converted into BLAST databases by the formatdb command in the NCBI BLAST package. We used the BLASTN program in the NCBI BLAST + version 2.2.26 package with the default parameters to search for similar sequences. The hits with an *E*-value under 10^{-4} were used for subsequent analyses.

LCA analysis for taxonomic classification

The LCA method assigns sequence reads to taxa with a criterion for the resolution of assignments [8]. $h(q, s)$ is the set of taxa found by a BLAST search for a sequence read q under the threshold of the bit score s . For a set of taxa $h(q, s)$, the common ancestor located farthest from the root of the taxonomic tree defines the LCA as the representative taxon. Thus, the LCA allows the assignment of a read to a single taxon. At the same time, the taxonomic levels indicate the resolution of assignments because the LCA allows broad hits to be assigned as high-level taxa but specific hits to be assigned as low-level taxa. It also means that the number of the LCA assigned reads depends on the thresholds of the bit scores for BLAST hits.

We use MEGAN software version 4.62.5 for the LCA analysis [8]. MEGAN assigns sequence reads into taxa at



different hierarchical levels such as family, genus, and species in the taxonomic ordering relation.

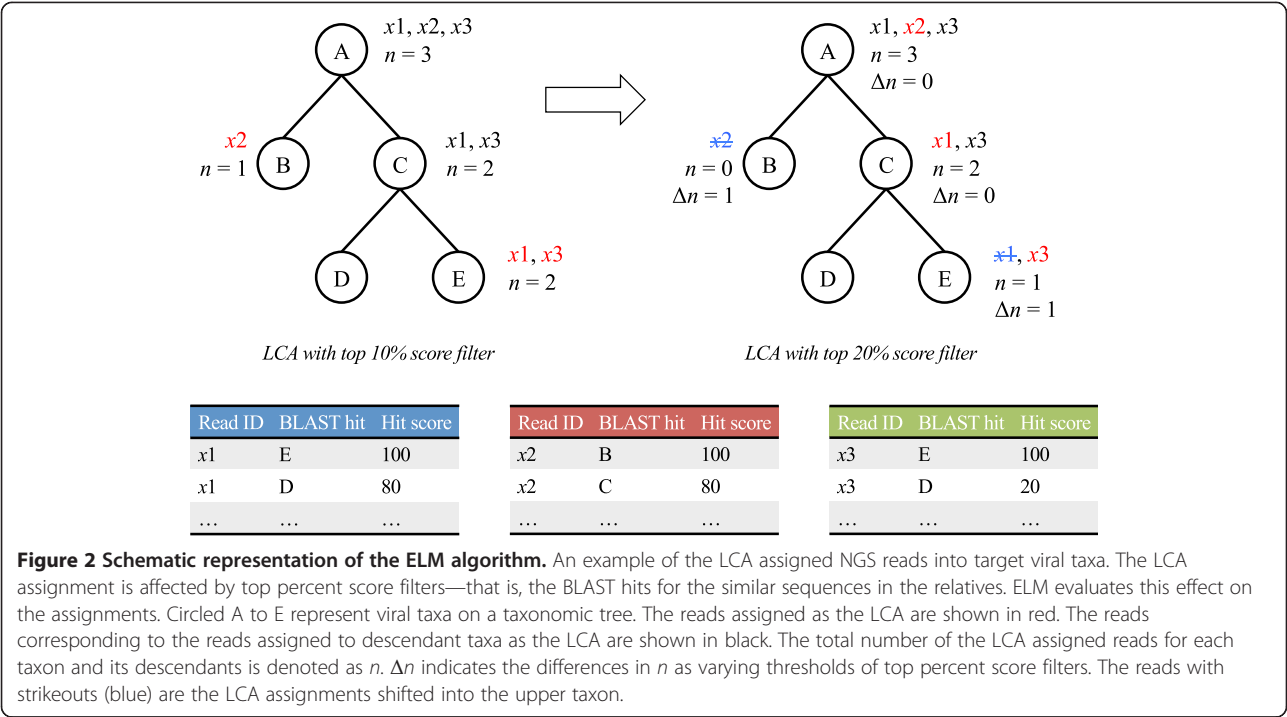
ELM for evaluating the LCA assignments

Since the LCA method solely provides the taxonomic levels of the assignments and the BLAST scores for short reads depend on local genomic regions, it is difficult to discriminate between true and false assignments with the LCA assignments or the BLAST scores alone. Here we combine the LCA assignments with top percent score filters for the BLAST hits to find out not only species-specific but also conserved reads, and introduce an additional criterion for the taxonomic assignment. First, ELM repeats the LCA analysis further under different top percent score filters for the BLAST hits and compares these LCA assignments with the reference assignment under the top 10% score filter (Figure 2).

Here, the top x percent score filter retains the BLAST hits whose bit scores lie within $x\%$ of the best score [8]. $n(x)$ is the total number of the LCA assigned reads for a taxon and its descendants under top x percent score filter. Then the difference Δn from that under the reference top 10% score filter is given by:

$$\Delta n = n(10) - n(x) \quad (1)$$

Here, Δn indicates to what extent the assigned reads are shifted into upper taxa as increasing x greater than 10%. We analyzed the increase of Δn , which is associated with sequence similarity to relatives, to discriminate between true and false assignments. In the statistical analysis of Δn , we introduce the inflation index IF , which is the Z score for outlier detection, to compare the effect



of top percent score filters on the taxonomic assignments. The IF for a taxon is given by:

$$IF = \frac{\Delta n - \mu}{\sigma} \quad (2)$$

where μ is the average of Δn for all assigned taxa, and σ is the standard deviation. Since multiple comparisons in IF s under top percent score filters ranging from 20% to 100% are performed nine times at 10% intervals, a P value of less than 0.05/9 is accepted for statistical significance after Bonferroni correction. Accordingly, $IF > 2.54$ (one-tailed) is accepted with statistical significance.

Utility and discussion

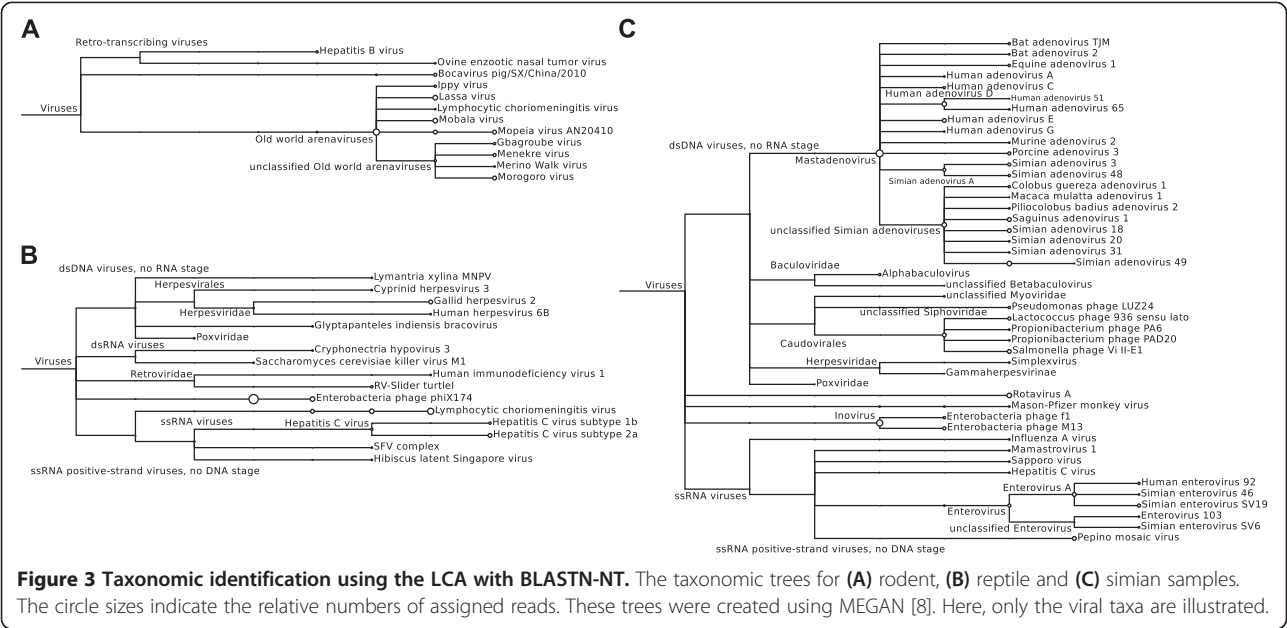
Benchmark tests for NGS datasets

To evaluate the ability of ELM to detect pathogenic viruses from large sequence datasets, five real datasets were used. Dataset 1 consisted of 4,449,766 unassembled reads from a rodent sample in Zambia [15]. Reads with an average length of 236 bases were obtained by Ion Torrent Personal Genome Machine (PGM) sequencing. Dataset 2 consisted of 4,146,547 unassembled reads from a reptile sample (SRR: 527074) deposited in the NCBI Sequence Read Archive (SRA). Reads with an average length of 200 bases were obtained by Illumina sequencing [16]. Dataset 3 consisted of 12,393,506 unassembled reads from a simian sample (SRR: 167721) deposited in the SRA. Reads with an average length of 73 bases were obtained by Illumina sequencing [17]. We selected these

three datasets to evaluate the effects of the read length, host and NGS platform. Furthermore, we applied ELM to fecal samples including multiple virus and phage taxa in dataset 4 (SRR: 1055974 for 12-day-old piglets) and dataset 5 (SRR: 1055972 for 54-day-old piglets). Reads with an average length of 291 bases in dataset 4 and 400 bases in dataset 5 were obtained by 454 GS FLX Titanium sequencing [18]. In these benchmark tests, the BLAST searches were performed on a workstation with an Intel Sandy Bridge CPU 2.6 GHz processor. We compared the result of the BLASTN search for the customized database with that for the NCBI NT database.

Identification of infecting viruses using the LCA with BLASTN-NT

To identify infecting viruses, we performed conventional LCA-based assignment using the results of a BLASTN search of the NCBI NT database (Figure 3). The taxa assigned at the 6th taxonomic level from the root in dataset 1 showed that this rodent host was infected with *Old world arenaviruses* (Figure 3A). A previous study showed that the rodent host was infected with *Luna virus*, which belongs to the *Old world arenaviruses* [15]. Totally, 99.9% of the sequences were derived from eukaryotes, including sequences from the rodent host. The reptile host in dataset 2 was infected with *Lymphocytic choriomeningitis virus* (Figure 3B). This result was consistent with the closest virus described in the literature [16]. In dataset 2, 99.5% of the sequences were probably derived from the reptile host. According to the literature concerning dataset 3, the



simian host was infected with a novel simian adenovirus, which is close to *Simian adenovirus 3*, *Simian adenovirus 18* and *Simian adenovirus 21* with about 55% pairwise nucleotide identity [17]. We found *Simian adenovirus 49*, *Simian adenovirus 18* and *Simian adenovirus 1* in dataset 3 (Figure 3C), suggesting results similar to those in the literature. Similarly, in dataset 3, most of the sequences (95.1%) were likely derived from the simian host.

To assess the required computational resources, we measured the elapsed time for the BLAST search (Table 1). As seen in Table 1, we found that the elapsed time for the BLAST search depended on the number of reads and hits. Although multiple threads and parallel jobs reduced the computational time, we needed at least one day with 8 threads and 32 parallel jobs. The sizes of the resulting tabulated format files ranged from 60–648 gigabytes, possibly affecting the elapsed time for the LCA analysis.

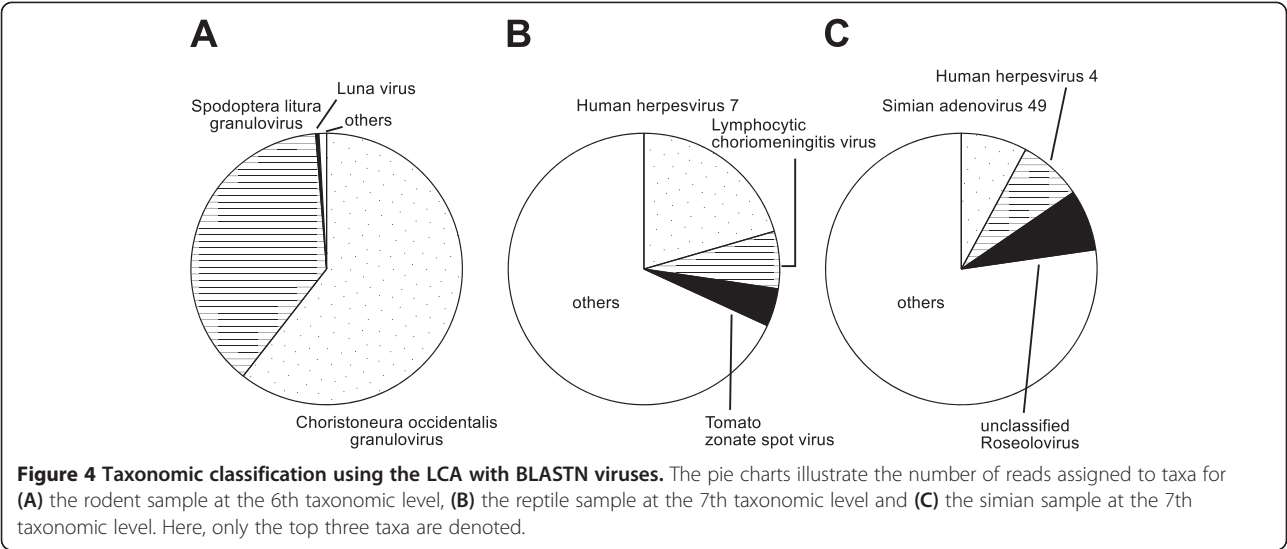
Taxonomic classification using the LCA with BLASTN viruses
According to the literature on taxonomic classification, the sequence similarity search of BLAST is a computational bottleneck [14]. Therefore, we used the customized viral database for a BLAST search to investigate how much the computational time was reduced and whether the conventional LCA could identify the infecting viruses.

Table 1 Elapsed time for the LCA with BLASTN-NT

Dataset No.	# of reads	# of BLASTN hits	CPU time	
			BLAST	LCA
1	4,449,766	4,424,602	12,179 h	96 m
2	4,146,547	2,754,210	8,704 h	22 m
3	12,393,506	10,674,129	23,313 h	82 m

In Figure 4, the top 3 assigned reads show the capturing of infecting viruses. However, most of the assigned taxa were false positives (Figure 4). At the 6th taxonomic level of assignments in dataset 1, 98.7% of the reads were assigned into *Choristoneura occidentalis granulovirus* and *Spodoptera litura granulovirus*, and only 0.4% (1,518/383,939) of them were assigned into *Luna virus* (Figure 4A). In the case of BLAST-NT, we failed to identify *Luna virus* but detected *Old world arenaviruses*, with 1,245 reads at the 5th taxonomic level, including the following relatives: *Mobala virus*, 125 reads; *Morogoro virus*, 73 reads; and *Mopeia virus*, 56 reads. In dataset 2, 8,387 reads were assigned into 141 viral taxa at the 7th taxonomic level, and 573 were assigned into *Lymphocytic choriomeningitis virus* (Figure 4B). In the case of BLAST-NT, 454 reads were assigned into *Lymphocytic choriomeningitis virus*. These results showed consistency between BLAST viruses and BLAST-NT. Of the 5,952 reads assigned into viral taxa at the 7th taxonomic level in dataset 3, 468 were assigned into *Simian adenovirus 49* (Figure 4C). The most assigned taxon in BLAST viruses was *Simian adenovirus 49*, but 99 reads in BLAST-NT were assigned into the closest relative, *Simian adenovirus 18*. Although the assignments with BLAST-NT were more favorable than those with BLAST viruses, the coverage of the identified *Simian adenovirus 49* was sufficient to perform the subsequent analysis. These results showed that the sensitivity of the LCA with BLASTN viruses outperformed the LCA with BLASTN-NT, suggesting that the BLAST search of the viral database was sufficient for subsequent analysis.

Next, we investigated whether the elapsed time for the BLASTN search was effectively reduced (Table 2). The elapsed time in the BLAST search was reduced to



0.03%-0.05% (Tables 1 and 2). This showed the synergy effect of the reduction of the custom database to 0.1% (from 38 Gb to 49 Mb) for the size of FASTA files and to 1.4%-8.8% for the number of BLASTN hits (Tables 1 and 2). Furthermore, the elapsed time for the LCA analysis was also reduced despite the additional nine assignments for ELM analysis.

Identification of infecting viruses using ELM with BLASTN viruses

To reduce the false assignments of the BLAST search for the customized viral database, we compared true and false assignments to confirm whether the true assignments altered into high-level taxa in ELM analysis (Figure 5). As shown for the 6th taxonomic level assignments of dataset 1 in Figure 5A and D, the assignment of *Luna virus* was significantly changed ($IF > 2.54$, ranging from 20% to 100%), suggesting that ELM correctly identified the infecting virus. In the 7th taxonomic level assignments of dataset 2, the assignment of *Lymphocytic choriomeningitis virus* was most changed ($IF > 9$, ranging from 20% to 100%) but, at the 6th taxonomic level, those of *unclassified Tospovirus*, *Tomato spotted wilt virus* and *Impatiens necrotic spot virus* were only slightly changed (Figure 5B and E). Figure 5C and F show that, in the 7th taxonomic level assignments of dataset 3, the assignment of *Simian adenovirus 49* was

significantly changed ($IF > 10$, ranging from 20% to 100%). However, at the 6th taxonomic level, the assignments of *Ictalurid herpesvirus 1*, *Simian adenovirus 3* and *Human adenovirus 54* were also changed, suggesting that ELM failed to exclude the false assignment of *Ictalurid herpesvirus 1*. On the other hand, ELM excluded the most of the false assignments ($IF < 2.54$, ranging from 20% to 100%). These results suggested that we could find out true assignments mainly by high inflation indices and partly by low taxonomic levels. The results for viruses identified using ELM with BLASTN viruses were assembled using SSAKE v3.8.1 [19] and are summarized in Table 3.

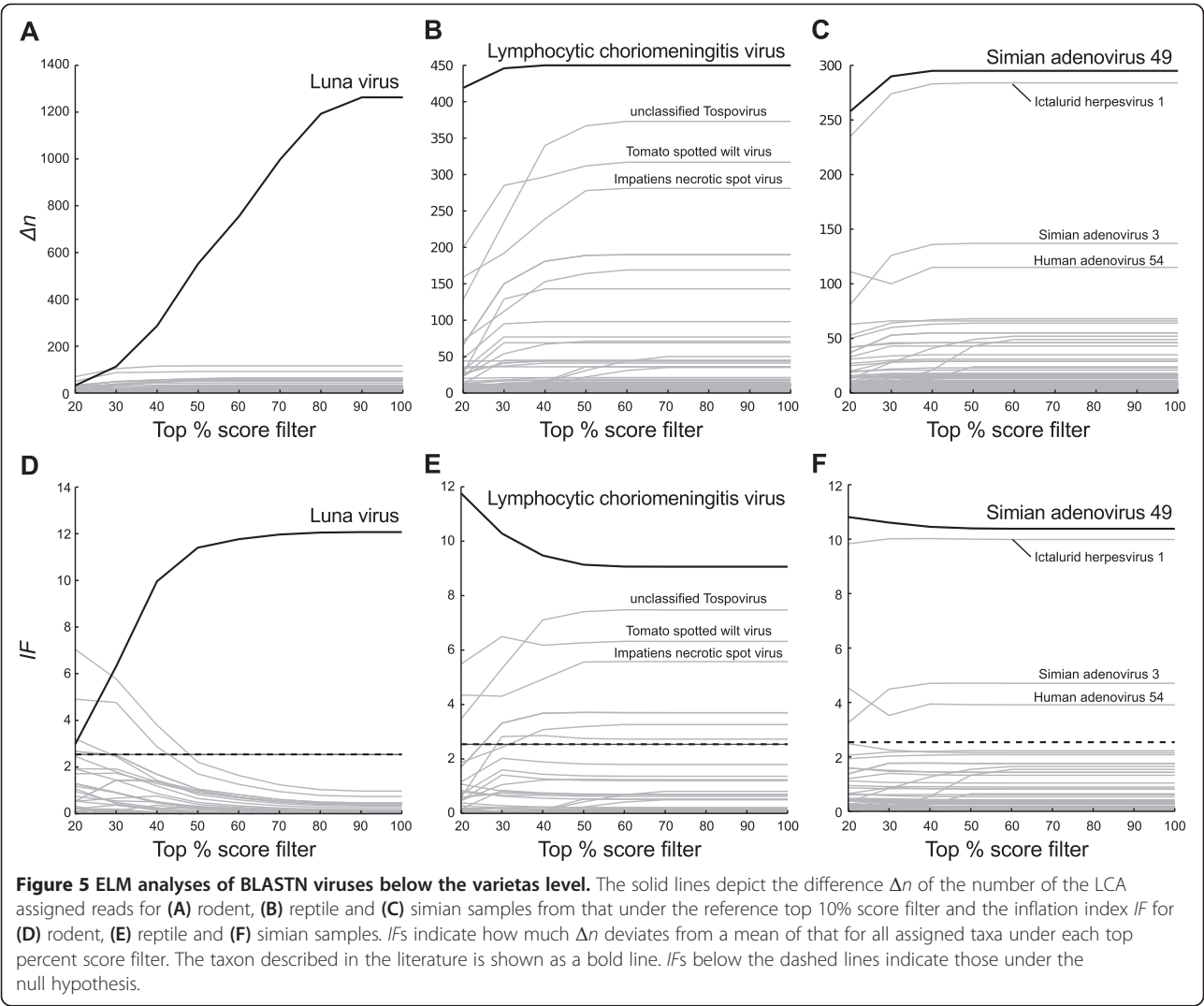
Next, we evaluated the effect of the BLAST hit score on the inflation indices. The results showed that the inflation indices had little association with the *E*-value in the BLAST search (Additional file 1: Figure S1). We also investigated the coverage of the BLAST hits. Valid hits in dataset 1 were distributed across target genomic sequences but not a specific genomic sequence, something not seen in datasets 2 and 3 (Additional file 1: Figure S2).

Virome analyses using ELM with BLASTN viruses

To investigate whether ELM could detect multiple viral taxa, we analyzed the fecal virome of piglets using ELM with BLASTN viruses. We identified the shift of *Kobuvirus* in dataset 4 to *Bocavirus* and *Dependovirus* in dataset 5, which depended on the age of the piglets (Table 4). These results were consistent with abundant virus genera described in the literature [18]. However, we failed to identify *pig stool-associated small circular DNA virus* in dataset 5 (Table 4). This virus belongs to the *single-stranded circular DNA viruses*. The members of this family show extensive genetic diversity [20]. The

Table 2 Elapsed time for ELM with BLASTN viruses

Dataset No.	# of BLASTN hits	CPU time	
		BLAST	LCA
1	387,271	4 h	5 m
2	38,948	4 h	1 m
3	379,780	6 h	5 m



results suggested that, in this case, the inflation index was not preferable for evaluating the LCA assignments.

Interpretation of results

ELM with a specific database drastically reduced the computational time and saved disk space. Furthermore, ELM was effective even for short reads. Though short

reads can reduce the accuracy of BLAST searches, in this study we verified ELM for average lengths of between 73 and 400 bases. The results showed no difference between the capabilities for taxonomic assignment.

One approach to reduce the computational time needed for the BLAST search is the subtraction of reads by mapping host-derived reads onto reference sequences

Table 3 Detection of the viral genomes using ELM with BLASTN viruses

Dataset No.	Virus ^a	# of reads	# of contigs	Average contig length
1	Luna virus	1,518	405	454 nt
2	LCMV	573	33	117 nt
3	SAdV-49	468	11	89 nt

Contigs were assembled using SSAKE v3.8.1 [19]. ^aLCMV, Lymphocytic choriomeningitis virus; SAdV-49, Simian adenovirus 49.

Table 4 Detection of abundant virus genera in fecal viromes of piglets using ELM with BLASTN viruses

Dataset No.	ELM with BLASTN viruses (# of reads)	LCA with BLAST-NT (# of reads)
4	Kobuvirus ^a (6,449)	Kobuvirus (6,446)
5	Dependovirus ^a (759), Bocavirus (133)	Dependovirus (754), Bocavirus (528), Chimpanzee stool associated circular ssDNA virus ^b (106)

^aThe genus includes descendant taxa $IF > 2.54$. ^bAccording to the literature [18], this virus is the novel pig stool-associated single-stranded DNA virus, which is not assigned to a specific genus.

[21,22]. This approach might be considered effective for reducing analyzed sequence data but is limited to known hosts. It is not suitable for surveillance of wild animals or metagenomic analysis because the host sequences have yet to be deposited in databases. Therefore, we need to decide a moderate threshold for NGS data before the mapping.

For ELM we adopted another approach using specific databases composed only of target sequences to reduce the computational cost. The difficulty in applying this approach directly to virus identification was the increase of false positive assignments (Figure 4). We tested several ways to solve this problem. As shown in Additional file 1: Figure S1, changing the threshold of the *E*-value dependent on the size of the database is probably not effective for discriminating between true and false assignments. The criteria for evaluating breadth coverage, i.e. the proportion of reads mapped across the hit genomes and the depth coverage (the number of reads mapped at a position), also failed to identify the target viruses (Additional file 1: Figure S2). On the other hand, ELM analyzed how sequence similarity to the relatives changes. This extension of the LCA method suppressed the rise in false positive assignments. A limitation of ELM would be the false-negative errors because ELM cannot detect viruses distantly related to other relatives (Table 4). Therefore, viruses without relatives should be carefully handled without the inflation index.

Conclusions

ELM is especially useful for the first screening of infectious diseases caused by viruses. In surveillance for pathogenic viruses, taxonomic assignment of the host sequence is not necessary for the initial screening. For this, sensitivity for detecting viruses is particularly required. Our results suggest that ELM recovers most reads assigned to target viruses. Therefore, we can apply these results to further sophisticated analyses. ELM will contribute to analyses of NGS data for limited targets such as the direct diagnosis of viral infections.

Availability and requirements

ELM is freely accessible at <http://bioinformatics.czc.hokudai.ac.jp/ELM/>. The web interface has been tested in the following web browsers: Google Chrome (version 36), Microsoft Internet Explorer (version 11) and Mozilla Firefox (version 31).

Additional file

Additional file 1: Figure S1. Distribution of *E*-values in BLASTN viruses. The figure shows the effect of *E*-value on discrimination between true and false assignments. **Figure S2.** Read coverage on the position of target genomes. The figure shows the effect of breadth and depth coverage on discrimination between true and false assignments.

Abbreviations

NGS: Next generation sequencing; LCA: Lowest common ancestor; ELM: Enhanced LCA-based method.

Competing interests

The authors declare that they have no competing interests.

Authors' contributions

KU designed, implemented, tested and evaluated the method, and wrote the manuscript. AI provided the experimental NGS data analyzed in the manuscript. KI participated in the design and evaluation of the method and collaborated in writing the manuscript. All authors read and approved the manuscript.

Acknowledgements

The work was supported by the Japan Initiative for Global Research Network on Infectious Diseases.

Author details

¹Division of Bioinformatics, Research Center for Zoonosis Control, Hokkaido University, Sapporo, Hokkaido 001-0020, Japan. ²Hokudai Center for Zoonosis Control in Zambia, Research Center for Zoonosis Control, Hokkaido University, Sapporo, Hokkaido 001-0020, Japan.

Received: 9 August 2013 Accepted: 7 July 2014

Published: 28 July 2014

References

- Garten RJ, Davis CT, Russell CA, Shu B, Lindstrom S, Balish A, Sessions WM, Xu X, Skepner E, Deyde V, Okomo-Adhiambo M, Gubareva L, Barnes J, Smith CB, Emery SL, Hillman MJ, Rivaller P, Smagala J, de Graaf M, Burke DF, Fouchier RA, Pappas C, Alpuche-Aranda CM, Lopez-Gatell H, Olivera H, Lopez I, Myers CA, Faix D, Blair PJ, Yu C, et al: **Antigenic and genetic characteristics of swine-origin 2009 A(H1N1) influenza viruses circulating in humans.** *Science* 2009, **325**(5937):197–201.
- Gao R, Cao B, Hu Y, Feng Z, Wang D, Hu W, Chen J, Jie Z, Qiu H, Xu K, Xu X, Lu H, Zhu W, Gao Z, Xiang N, Shen Y, He Z, Gu Y, Zhang Z, Yang Y, Zhao X, Zhou L, Li X, Zou S, Zhang Y, Yang L, Guo J, Dong J, Li Q, Dong L, et al: **Human infection with a novel avian-origin influenza A (H7N9) virus.** *N Engl J Med* 2013, **368**(20):1888–1897.
- Li W, Shi Z, Yu M, Ren W, Smith C, Epstein JH, Wang H, Cramer G, Hu Z, Zhang H, Zhang J, McEachern J, Field H, Daszak P, Eaton BT, Zhang S, Wang LF: **Bats are natural reservoirs of SARS-like coronaviruses.** *Science* 2005, **310**(5748):676–679.
- Feldmann H, Wahl-Jensen V, Jones SM, Stroher U: **Ebola virus ecology: a continuing mystery.** *Trends Microbiol* 2004, **12**(10):433–437.
- Nash D, Mostashari F, Fine A, Miller J, O'Leary D, Murray K, Huang A, Rosenberg A, Greenberg A, Sherman M, Wong S, Layton M: **The outbreak of West Nile virus infection in the New York City area in 1999.** *N Engl J Med* 2001, **344**(24):1807–1814.
- Yu XJ, Liang MF, Zhang SY, Liu Y, Li JD, Sun YL, Zhang L, Zhang QF, Popov VL, Li C, Qu J, Li Q, Zhang YP, Hai R, Wu W, Wang Q, Zhan FX, Wang XJ, Kan B, Wang SW, Wan KL, Jing HQ, Lu JX, Yin WW, Zhou H, Guan XH, Liu JF, Bi ZQ, Liu GH, Ren J: **Fever with thrombocytopenia associated with a novel bunyavirus in China.** *N Engl J Med* 2011, **364**(16):1523–1532.
- Barzon L, Lavezzo E, Militello V, Toppo S, Palu G: **Applications of next-generation sequencing technologies to diagnostic virology.** *Int J Mol Sci* 2011, **12**(11):7861–7884.
- Huson DH, Auch AF, Qi J, Schuster SC: **MEGAN analysis of metagenomic data.** *Genome Res* 2007, **17**(3):377–386.
- Turnbaugh PJ, Ley RE, Mahowald MA, Magrini V, Mardis ER, Gordon JL: **An obesity-associated gut microbiome with increased capacity for energy harvest.** *Nature* 2006, **444**(7122):1027–1031.
- Handley SA, Thackray LB, Zhao G, Presti R, Miller AD, Droit L, Abbink P, Maxfield LF, Kambal A, Duan E, Stanley K, Kramer J, Macri SC, Permar SR, Schmitz JE, Mansfield K, Brenchley JM, Veazey RS, Stappenbeck TS, Wang D, Barouch DH, Virgin HW: **Pathogenic simian immunodeficiency virus infection is associated with expansion of the enteric virome.** *Cell* 2012, **151**(2):253–266.
- Tong S, Li Y, Rivaller P, Conrardy C, Castillo DA, Chen LM, Recuenco S, Ellison JA, Davis CT, York IA, Turmelle AS, Moran D, Rogers S, Shi M, Tao Y, Weil MR, Tang K, Rowe LA, Sammons S, Xu X, Frace M, Lindblade KA, Cox NJ, Anderson LJ, Rupprecht CE, Donis RO: **A distinct lineage of influenza A virus from bats.** *Proc Natl Acad Sci U S A* 2012, **109**(11):4269–4274.

12. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ: **Basic local alignment search tool.** *J Mol Biol* 1990, **215**(3):403–410.
13. Monzoorul Haque M, Ghosh TS, Komanduri D, Mande SS: **Sort-ITEMS: Sequence orthology based approach for improved taxonomic estimation of metagenomic sequences.** *Bioinformatics* 2009, **25**(14):1722–1730.
14. Gerlach W, Stoye J: **Taxonomic classification of metagenomic shotgun sequences with CARMA3.** *Nucleic Acids Res* 2011, **39**(14):e91.
15. Ishii A, Thomas Y, Moonga L, Nakamura I, Ohnuma A, Hang'ombe B, Takada A, Mweene A, Sawa H: **Novel arenavirus, Zambia.** *Emerg Infect Dis* 2011, **17**(10):1921–1924.
16. Stenglein MD, Sanders C, Kistler AL, Ruby JG, Franco JY, Reavill DR, Dunker F, Derisi JL: **Identification, characterization, and in vitro culture of highly divergent arenaviruses from boa constrictors and annulated tree boas: candidate etiological agents for snake inclusion body disease.** *MBio* 2012, **3**(4):e00180–00112.
17. Chen EC, Yagi S, Kelly KR, Mendoza SP, Tarara RP, Canfield DR, Maninger N, Rosenthal A, Spinner A, Bales KL, Schnurr DP, Lerche NW, Chiu CY: **Cross-species transmission of a novel adenovirus associated with a fulminant pneumonia outbreak in a new world monkey colony.** *PLoS Pathog* 2011, **7**(7):e1002155.
18. Sachsenroder J, Twardziok SO, Scheuch M, Johne R: **The general composition of the faecal virome of pigs depends on age, but not on feeding with a probiotic bacterium.** *PLoS One* 2014, **9**(2):e88888.
19. Warren RL, Sutton GG, Jones SJ, Holt RA: **Assembling millions of short DNA sequences using SSAKE.** *Bioinformatics* 2007, **23**(4):500–501.
20. Cheung AK, Ng TF, Lager KM, Bayles DO, Alt DP, Delwart EL, Pogranichniy RM, Kehrli ME Jr: **A divergent clade of circular single-stranded DNA viruses from pig feces.** *Arch Virol* 2013, **158**(10):2157–2162.
21. Chang Y, Cesarman E, Pessin MS, Lee F, Culpepper J, Knowles DM, Moore PS: **Identification of herpesvirus-like DNA sequences in AIDS-associated Kaposi's sarcoma.** *Science* 1994, **266**(5192):1865–1869.
22. Simons JN, Pilot-Matias TJ, Leary TP, Dawson GJ, Desai SM, Schlauder GG, Muerhoff AS, Erker JC, Buijk SL, Chalmers ML, van Sant CL, Mushahwar IK: **Identification of two flavivirus-like genomes in the GB hepatitis agent.** *Proc Natl Acad Sci U S A* 1995, **92**(8):3401–3405.

doi:10.1186/1471-2105-15-254

Cite this article as: Ueno et al.: ELM: enhanced lowest common ancestor based method for detecting a pathogenic virus from a large sequence dataset. *BMC Bioinformatics* 2014 **15**:254.

Submit your next manuscript to BioMed Central and take full advantage of:

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at
www.biomedcentral.com/submit

