

ARPIP: Ancestral Sequence Reconstruction with Insertions and Deletions under the Poisson Indel Process

GHOLAMHOSSEIN JOWKAR^{1,2,3,*}, JŪLIJA PEČERSKA^{1,2}, MASSIMO MAIOLO^{1,2,4}, MANUEL GIL^{1,2}, AND MARIA ANISIMOVA^{1,2}

¹Zurich University of Applied Sciences, School of Life Sciences and Facility Management, CH-8820 Wädenswil, Switzerland; ²Swiss Institute of Bioinformatics, CH-1015 Lausanne, Switzerland; ³University of Neuchâtel, Institute of Biology, CH-2000 Neuchâtel, Switzerland; and ⁴University of Bern, Institute of Pathology, CH-3008 Bern, Switzerland

*Correspondence to be sent to: ZHAW, School of Life Sciences and Facility Management, Applied Computational Genomics Group, Schloss 1, CH-8820 Wädenswil, Switzerland;
E-mail: jowk@zhaw.ch.

Received 4 November 2021; accepted 6 July 2022

Associate Editor: Adrian Paterson

Abstract.—Modern phylogenetic methods allow inference of ancestral molecular sequences given an alignment and phylogeny relating present-day sequences. This provides insight into the evolutionary history of molecules, helping to understand gene function and to study biological processes such as adaptation and convergent evolution across a variety of applications. Here, we propose a dynamic programming algorithm for fast joint likelihood-based reconstruction of ancestral sequences under the Poisson Indel Process (PIP). Unlike previous approaches, our method, named ARPIP, enables the reconstruction with insertions and deletions based on an explicit indel model. Consequently, inferred indel events have an explicit biological interpretation. Likelihood computation is achieved in linear time with respect to the number of sequences. Our method consists of two steps, namely finding the most probable indel points and reconstructing ancestral sequences. First, we find the most likely indel points and prune the phylogeny to reflect the insertion and deletion events per site. Second, we infer the ancestral states on the pruned subtree in a manner similar to FastML. We applied ARPIP (Ancestral Reconstruction under PIP) on simulated data sets and on real data from the *Betacoronavirus* genus. ARPIP reconstructs both the indel events and substitutions with a high degree of accuracy. Our method fares well when compared to established state-of-the-art methods such as FastML and PAML. Moreover, the method can be extended to explore both optimal and suboptimal reconstructions, include rate heterogeneity through time and more. We believe it will expand the range of novel applications of ancestral sequence reconstruction. [Ancestral sequences; dynamic programming; evolutionary stochastic process; indel; joint ancestral sequence reconstruction; maximum likelihood; Poisson Indel Process; phylogeny; SARS-CoV.]

Phylogenetics is a wide research field with a variety of applications ranging from reconstructing the tree of life to investigating ongoing epidemics. Phylogenetic trees provide insight into unobservable evolutionary events in the past such as adaptation or mass extinction events. Phylogenetic inference can be divided into several inter-related tasks including sequence alignment, phylogeny estimation, detection of selection, and ancestral sequence reconstruction (ASR). ASR aims to infer the likely ancestral sequences for a set of existing homologous sequences.

ASR allows researchers to pursue a wide range of topics from determining the origins of life or epidemics to developing personalized medicine (Pagel 1999; Liberles 2007). For example, the functionality of ancient genes can be investigated by reconstructing and synthesizing the genetic material inferred by ASR (Thornton 2004). Such analyses can help us understand the mechanisms underlying adaptation and speciation processes, inspiring new approaches for protein engineering (Chang et al. 2005) and drug design (Zakas et al. 2017). ASR can be used to study epidemiological origins of pathogens, particularly in light of recent coronavirus pandemics (Pagel 1999; Brintnell et al. 2021; Starr et al. 2022).

State-of-the-art likelihood-based ASR methods use Markov processes to model character substitutions through time. Such models account for various biases in character substitution, as well as divergence represented by evolutionary time (Yang et al. 1995; Pupko et al.

2000; Yang 2007). However, Markov models of molecular evolution do not include insertions or deletions (indels) as part of the evolutionary process, meaning that methods relying on these models have to treat gap characters separately. Most of the ASR methods adopt one of two preprocessing approaches. They either treat gaps as missing/ambiguous data or remove gap characters entirely. However, indels represented by gaps carry an important evolutionary signal (Dessimoz and Gil 2010) and are in fact a major driving force of genomic divergence (Tao et al. 2007). Therefore, methods that model indels explicitly have a clear advantage over methods that do not. Up to this point, most existing frequentist algorithms do not include indel modeling except for two methods, Ancestors (Diallo et al. 2009) and FastML (Ashkenazy et al. 2012). Ancestors have exponential computational complexity and therefore have not been widely adopted by users. FastML handles indels using a heuristic approach called indel coding. The method relies on the linear time complexity algorithm (Pupko et al. 2000) for joint maximum likelihood (ML) ASR using dynamic programming (DP). FastML makes the analyses of large data sets tractable. Currently, it is provided as a web service (Ashkenazy et al. 2012). While the results of indel coding can be interpreted from an evolutionary standpoint retrospectively, the approach does not, however, include an explicit evolutionary indel model. All things considered, most methods rely on standard models of sequence evolution

without indels which is an issue that can only be resolved by including character and indel evolution in a single model.

Two pioneering mathematical models describing the evolution of indels are TKF91 and TKF92 (Thorne et al. 1991, 1992). However, the computation of marginal likelihood under these models has exponential time complexity, rendering the methods relying on these models extremely computationally intensive, and making inference under these models unrealistic on large data sets. More recently, (Bouchard-Côté and Jordan, 2013) proposed the Poisson Indel Process (PIP) model which is based on TKF91. PIP describes insertions by a Poisson process defined on the tree topology, while substitutions and deletions are described by a continuous-time Markov process where deletions are modeled as an absorbing state. The assumption of independence between insertion and substitution/deletion enabled a major computational improvement over the previous models (Bouchard-Côté and Jordan 2013). In PIP, the insertion rate is also independent of the length of a sequence, which is a realistic assumption based on the data that is most commonly analyzed (Bouchard-Côté 2010, p. 93). In contrast to TKF91, the PIP model allows to compute marginal likelihoods in linear time with respect to the number of sequences, which enables a variety of phylogenetic applications (e.g., Maiolo et al. 2018).

In this study, we use the PIP model for joint reconstruction of ancestral character states including insertions and deletions. Our method ARPIP (Ancestral Reconstruction under PIP) is implemented in the ML framework, that is, we use an empirical Bayesian approach with ML estimates. Given a multiple sequence alignment (MSA) and a phylogenetic tree, we first use PIP to infer insertion and deletion points on the tree. Insertion and deletion points are the specific locations on the phylogeny where the events have happened. Next, we extract a subtree rooted at the insertion point and pruned by the deletion points. Finally, we reconstruct ancestral states on the extracted subtree using a modified version of Felsenstein's recursion (Felsenstein 1981), similar to the FastML algorithm (Pupko et al. 2000). In the following, we describe the method in detail, validate it by simulations, and demonstrate its performance in simulations and on a real data set.

MATERIALS AND METHODS

ARPIP consists of two main algorithms: indel point inference and ancestral character inference.

The IndelPoints algorithm infers the most likely indel points for each site m of the given alignment (Appendix S2 of the Supplementary material available on Dryad at <http://dx.doi.org/10.5061/dryad.wstqjq2nj>). It traverses the tree in postorder and evaluates a set of possible indel scenarios for each node in the tree. A particular indel scenario defines a homology path \mathcal{H} . A homology path contains a single insertion point and a number of deletion points consistent with

the input MSA. IndelPoints finds the most likely indel scenario by maximizing the probability of \mathcal{H} given m . The maximization is simplified by reducing the MSA to gap and nongap states, and ignoring the substitution history without changing the result of the computation. This allows us to avoid matrix exponentiation, which is computationally expensive but necessary for the full likelihood computation.

Similar to the recursive likelihood computation, we traverse the tree and evaluate all the possible indel scenarios, selecting the best one at each node. At the tree root, we select the best homology path over the whole tree based on the best paths selected in the child nodes. For each site m , we use the inferred homology path to extract a subtree τ_m rooted at the insertion point \mathcal{I} and pruned by deletion points \mathcal{D} , which represents the most likely indel history for the given site.

Next, we reconstruct ancestral characters on the pruned subtrees in a manner similar to FastML (Pupko et al. 2000). For each site m , we use DP to reconstruct ancestral characters in two phases. The first phase can be seen as a modification of Felsenstein's peeling recursion for computing marginal likelihoods (Felsenstein 1981). As in the peeling recursion algorithm, we traverse the tree τ_m in postorder, starting from the leaves upward to the root and propagate partial likelihoods. However, instead of marginalizing over internal character states, for each MSA column m , we store the likelihood values $L_{k,v}$ and the corresponding best ancestral character states CS_v for each node v . In the second phase, the algorithm traverses the tree in preorder and for each node selects the ancestral character A_v with the highest conditional probability.

Preliminaries: The PIP Model

The PIP model describes the evolutionary process of substitutions, insertions, and deletions along the branches of a phylogenetic tree τ . Here, we include the basic description of the process, additional information on the PIP likelihood is available in Appendix S1 of the Supplementary material available on Dryad and a detailed description of PIP can be found in (Bouchard-Côté and Jordan, 2013).

Let $\tau = (\mathcal{V}, \mathcal{E}, b)$ represent a rooted binary phylogenetic tree, where set \mathcal{V} is the set of all vertices of the tree, \mathcal{E} is the set of all tree branches ($\mathcal{V} \times \mathcal{V}$), and b refers to the branch lengths in units of time (measured in expected substitutions and deletions per site).

The observed sequences are strings of characters from an alphabet Σ , which can be nucleotides, amino acids, or codons. The N observed sequences at the leaves of τ are denoted by set $\mathcal{L} \subset \mathcal{V}$, whereas set $\mathcal{V} \setminus \mathcal{L}$ is the set of $N - 1$ internal vertices. The root, the most recent common ancestor of all leaves, is labeled by Ω . The branch length $b(v)$ associated with node $(v \in \mathcal{V})$ spans from v to its parent vertex $pa(v)$ (see Fig. 1).

PIP is parameterized by insertion rate λ and deletion rate μ , with the process running over tree topology τ . For every node $v \in \mathcal{V}$, the probability of inserting a single

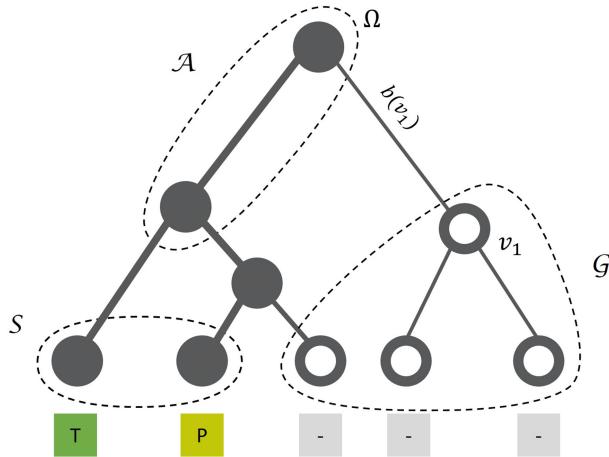


FIGURE 1. The phylogenetic tree τ rooted at Ω . $b(v_1)$ represents the branch length from Ω to v_1 . The leaves of the tree show a single column of the MSA including gaps as an additional character state. The set S is defined as all leaves with a character in the given column (not a gap). The set of potential insertion nodes \mathcal{A} contains the nodes ancestral to all nodes in S . Finally, the set of potential deletion nodes G is defined as all nodes which are either a leaf with a gap in the given column, or a node whose both children are in G .

character on edge $e = (\text{pa}(v) \rightarrow v)$ is proportional to the branch length and defined by $\iota(v)$ (see Appendix S1 of the Supplementary material available on Dryad). Similarly, the survival probability for a character inserted on edge e is $\beta(v)$ (see Appendix S1 of the Supplementary material available on Dryad). Additionally, we define the pure survival probability $\zeta = \exp(-\mu b(v))$ associated with node v as if the character was already present at the parent node $\text{pa}(v)$ (Maiolo 2019). Point substitutions and deletions are modeled by a continuous-time Markov process on $\Sigma_\epsilon = \Sigma \cup \{\epsilon\}$, where ϵ denotes the gap symbol. The generator matrix Q could be any arbitrary reversible substitution model, for example, WAG for amino acids (Whelan and Goldman 2001), or K80 for nucleotide data (Kimura 1980). Accordingly, the extended generator matrix is denoted by Q_ϵ and the extended quasistationary distribution is $\pi_\epsilon = [\pi, 0]$ (Bouchard-Côté and Jordan 2013).

Let \mathcal{G} define the site-specific set of all potential deletion points on the tree. \mathcal{G} consists of all leaves with a gap at the respective site, and of all the internal nodes whose all descendant leaves have a gap at that site. Next, consider the subset S of leaves that have a nongap character, $S = \{v \in \mathcal{L} : m_v \neq \epsilon\}$. Given the set S , we define the set \mathcal{A} of potential insertion points to include all nodes that are ancestral to all the leaves in S (see Fig. 1). In general, we compute the probability $p(m)$ of each individual MSA column by marginalizing over all possible homology paths underlying that MSA column (see Appendix S1 of the Supplementary material available on Dryad) based on their homology path probabilities f_v .

Inferring the Indel Points

We propose a progressive algorithm to infer the most likely indel points (homology path) on the tree under the

PIP model. For each site, we progressively find the best partial homology path (constrained by a subtree) and build on the intermediate results to get the most likely indel history on the whole tree. Since we search only for the most likely homology path, we compute a simplified likelihood function which accounts only for insertions and deletions and ignores substitutions.

Under PIP, two mutually exclusive node sets exist on the tree topology τ : the set of nodes where the character has gone extinct and the set of nodes where the character has definitely survived. The first is the set of potential deletion nodes \mathcal{G} , defined in the previous section. The second contains all the remaining nodes in the tree $v \in \mathcal{V} \setminus \mathcal{G}$ ($v \notin \mathcal{G}$). A node $v \notin \mathcal{G}$ may also be a potential insertion location, that is, $v \in \mathcal{A}$ (see Fig. 2 and Appendix S2 of the Supplementary material available on Dryad for the detailed description). While the set \mathcal{A} may contain multiple nodes, a homology path can only have a single insertion location. Consequently, each node in \mathcal{A} is associated with a single homology path with the highest probability. This implies that when computing the probability of a homology path for node $v \in \mathcal{A}$, it is treated as the only potential insertion location, while all other nodes are treated as regular nodes in the tree. Notably, one cannot simply select the node with the highest insertion probability $\iota(v)$, as the probability of any given homology path also depends on the survival/extinction of the site in the children. Even though we separately describe the treatment of the two node types, all the necessary computation can be done in a single postorder traversal of the tree.

For each node v in the tree, we first compute f_v , the conditional probability of the deletion/substitution process, assuming that the character exists in v . We compute f_v for the most likely deletion scenario in this subtree rather than marginalize over all possible deletion locations. We also compute p_v , the conditional probability of the homology path assuming that the character was inserted at node v . A character necessarily has to be inserted at one of the nodes $v \in \mathcal{A}$, which means that the probability will be nonzero only for the potential insertion nodes.

In the progressive algorithm, we maintain several node sets that are needed to define the most likely homology path per node. Let \mathcal{I}_v denote the set of insertion points for the subtree rooted at v . Then, $\mathcal{I}_v = \emptyset$ for $v \notin \mathcal{A}$ and $\mathcal{I}_v = \{v\}$ for $v \in \mathcal{A}$. Similarly, \mathcal{D}_v denotes the set of deletion points for the subtree rooted at v .

For each node v , we store the locally optimal homology path $\mathcal{H}_v = \{\mathcal{I}_v, \mathcal{D}_v\}$ for the subtree rooted at v . Once we reach the root node Ω , all possible insertion locations would be considered, and the one with the highest probability is selected among those. At this point, the best homology path $\mathcal{H}_{\arg\max(p_v)}$ (defined by the highest conditional probability p_v) is used to extract the subtree τ_m rooted at $\mathcal{I}_{\arg\max(p_v)}$ and pruned by $\mathcal{D}_{\arg\max(p_v)}$, which represents the best possible indel points for the MSA column m . We will use τ_m to infer ancestral character states for column m . The IndelPoints algorithm

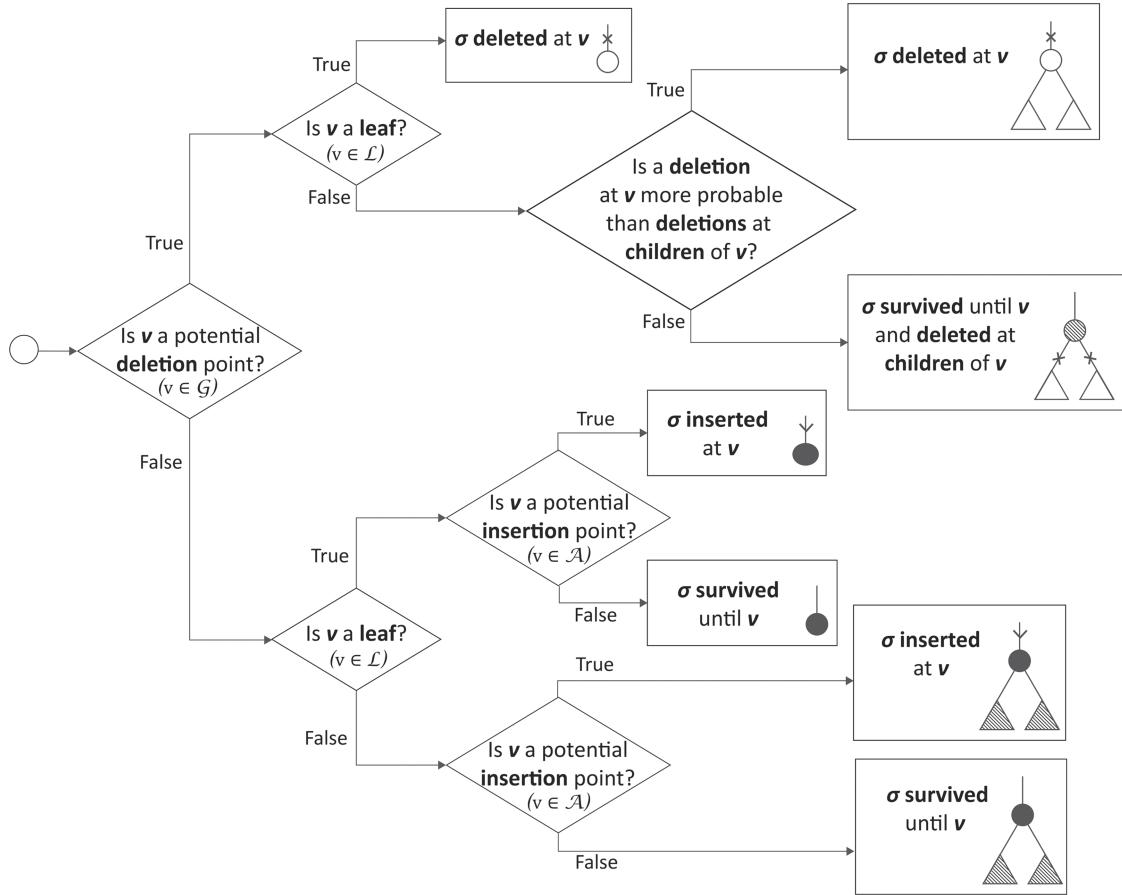


FIGURE 2. Overview of the IndelPoints algorithm. The tree is traversed in postorder to infer the most likely homology path progressively using the predefined sets: \mathcal{L} the set of all leaves, \mathcal{A} the set of potential insertion points, and \mathcal{G} the set of potential deletion points. Here, σ represents the character in focus and v is the node visited during the tree traversal.

is presented in Figure 2 and the pseudocode for the algorithm can be found in the [Appendix S2](#) of the [Supplementary material](#) available on Dryad.

DP Joint ASR

Our method performs ASR in a manner very similar to FastML ([Pupko et al. 2000](#)) with two crucial differences. First, we only work on a subtree τ_m of the original tree τ , which limits the reconstruction to the most probable insertion location at this site. This means we do not reconstruct any ancestral states where there were none. Second, to appropriately account for character deletion, the ancestral reconstruction is done using the PIP substitution rate matrix Q_ϵ .

The joint ASR method under PIP given column m and the pruned rooted phylogenetic subtree τ_m consists of two steps. The first step is to compute the partial likelihood values on subtree τ_m with the modified version of Felsenstein recursion algorithm, where both likelihood values and their corresponding ancestral character states are stored. The second step is to reconstruct the character states by picking the character with highest conditional

probability. The recursive algorithm for joint ASR is shown in [Appendix S3](#) of the [Supplementary material](#) available on Dryad together with the newly defined pseudocode for the procedure.

RESULTS

Three data sets were used to evaluate and illustrate our method. The first data set was simulated under the PIP. This data set allows to evaluate the performance of ARPIP under the true model. Given the true simulated trees and MSAs, both the homology path inference and ASR were evaluated.

The second data set was used to evaluate the performance of ARPIP for sequences with long indels. The data were generated by INDELible ([Fletcher and Yang 2009](#)) with two different settings using the same trees as for the PIP simulations. For this data set, the ancestral sequences are also known and can be used for evaluation. However, the PIP parameters for this data set must be inferred. As INDELible does not provide a comprehensive description of indel events on the phylogeny, we used these simulations to evaluate ancestral state inference.

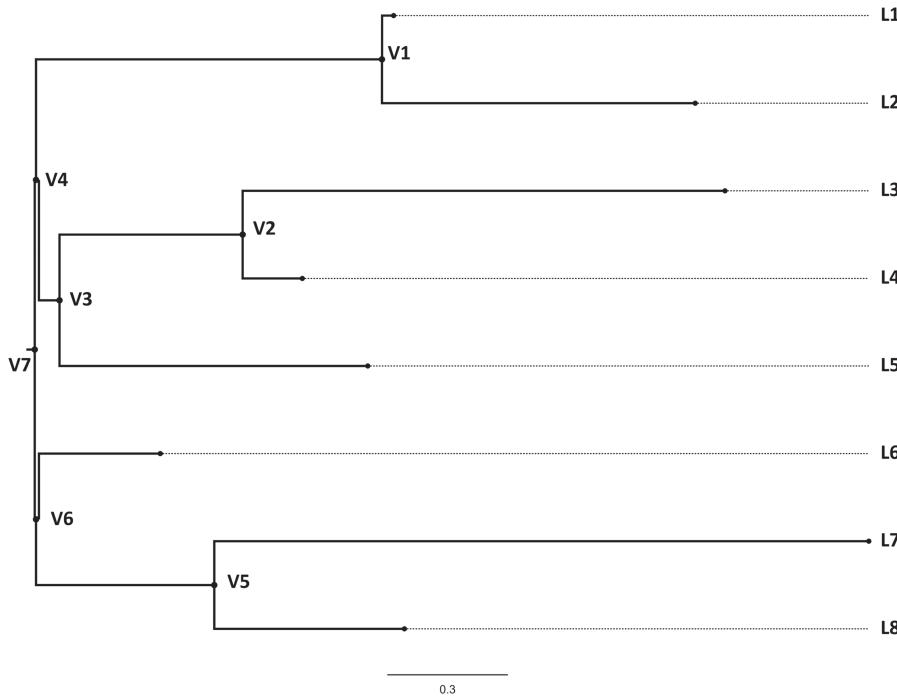


FIGURE 3. An example tree from the data set generated by the PIP simulator.

The third data set is a small coronavirus sample which was extracted from Uniprot (Bateman et al. 2020). With this data set, we aimed to provide a showcase of the method.

When analyzing both INDELible and real-life data, we first have to infer the PIP parameters, that is, λ and μ , given an MSA and a tree. This computation was done based on Brent's optimization method (Brent 1973), optimizing one parameter at a time until convergence. For all examples, the protein substitution model used is WAG (Whelan and Goldman 2001).

Note that ARPIP was developed for rooted trees. If an unrooted tree is provided, ARPIP uses midpoint rooting method to root the tree. Furthermore, ARPIP can perform ASR without a provided tree. The user can select from established fast methods like neighbor joining, BioNJ, UPGMA, and WPGMA to estimate the tree from the input MSA.

Data Simulated under PIP

The simulated sequences are given as input to ARPIP along with the true model parameters so that we only have to estimate the ancestral state values. The simulated data set contains 100 MSA/tree replicates with their corresponding evolutionary events. Each replicate was simulated using an eight taxa tree with a topology sampled from the uniform distribution and branch lengths sampled from an exponential distribution with the rate $\rho=2$, where ρ is a proxy for phylogenetic divergence. One of the simulated trees is shown in Figure 3. On average, the branch lengths of the simulated

trees were 0.45 units of time, ranging from minimal branch length of 0 and maximum branch length of 3.23. For the simulations, we set the deletion rate $\mu=0.1$ and the insertion rate $\lambda=10$ for PIP.

Analysis of the PIP simulated data set.—In order to assess the accuracy of ARPIP, we independently evaluated each inference step, IndelPoints and the joint ASR (see Table 1). To assess the accuracy of the IndelPoints algorithm, we also evaluated the inference of insertion and deletion events independently. As this data set was simulated under the same model we use for inference, we used the true parameter values in the analysis without inferring them ($\mu=0.1$ and $\lambda=10$). This way we can evaluate the method without the additional variation of parameter inference, which is done for the other two data sets.

Among the 100 input sets, 96.69% insertion points and 95.54% deletion points were inferred correctly. In the next step, we computed the accuracy of ASR per site. This number has been averaged over all existing sites over all MSA replicates. To evaluate the reconstructed ancestral sequences, we used three different metrics. Firstly, we counted the number of full ancestral columns

TABLE 1. ARPIP accuracy for inference on PIP simulated data

Metric	Accuracy (%)
Correctly inferred insertion points	96.08 ± 2.84
Correctly inferred deletion points	95.54 ± 2.80
Correctly inferred MSA columns	60.08 ± 9.60
Correctly inferred characters including gap	88.14 ± 3.91
Correctly inferred gap character	99.86 ± 0.26

TABLE 2. ARPIP accuracy for inference on INDELible simulated data

Metric	Accuracy (%)	
	Indel rate 0.01	Indel rate 0.05
Correctly inferred MSA columns	46.58 ± 13.22	59.49 ± 10.63
Correctly inferred characters including gap	83.49 ± 6.04	87.93 ± 4.33
Correctly inferred gap character	99.98 ± 0.16	99.95 ± 0.14

TABLE 3. *Betacoronavirus* sequences used in the analysis

Subgenus	Species	Uniprot accession number
<i>Embecovirus</i>	<i>Betacoronavirus 1</i> <i>China Rattus coronavirus HKU24</i> <i>Human coronavirus HKU1</i> <i>Murine coronavirus</i> <i>Myodes coronavirus 2JL14</i>	A0A191URB2 A0A0A7UZR7 U3NAI2 P11224 A0A2H4MXV6
<i>Hibcovirus</i>	<i>Bat Hp-betacoronavirus Zhejiang2013</i>	A0A088DJY6
<i>Merbecovirus</i>	<i>Hedgehog coronavirus 1</i> <i>Middle East respiratory syndrome-related coronavirus</i> <i>Pipistrellus bat coronavirus HKU5</i> <i>Tylonycteris bat coronavirus HKU4</i>	A0A4D6G1A4 K9N5Q8 A3EXD0 A3EX94
<i>Nobecovirus</i>	<i>Rousettus bat coronavirus GCCDC1</i> <i>Rousettus bat coronavirus HKU9</i>	A0A1B3Q5W5 A3EXG6
<i>Sarbecovirus</i>	<i>Severe acute respiratory syndrome-related coronavirus</i> <i>Severe acute respiratory syndrome-related coronavirus 2</i>	A0A3Q8AKM0 P0DTIC2

that were inferred correctly, which is 60.08% for this data set. Secondly, we counted the number of characters that were inferred correctly, which amounts to 88.14% of characters. Thirdly, we counted the number of gap characters themselves that were inferred correctly, which amounts to 99.86%.

Data Generated by INDELible

The data simulated by INDELible contain two sets of 100 MSA/tree replicas. Each replica was simulated using an eight taxa tree from PIP simulations. We used the Zipfian (power law) distribution for the indel model with $a=1.7$, to generate the samples where $a>1$ is the exponent characterizing the distribution. Empirical estimates of value a range from 1.5 to 2 (Fletcher and Yang 2009), which prompted us to select $a=1.7$. The maximum indel length was set to 5 to avoid MSAs with excessively long gaps. Two different indel rates of 0.01 and 0.05 were used for the simulation with INDELible.

Analysis of the INDELible simulated data set.—For this data set, ARPIP inferred the PIP parameters λ and μ as well as ancestral character states. For the two data sets of 100 MSA/tree replicates with indel rates of 0.01 and 0.05, ARPIP correctly inferred 46.58% and 59.49% of the ancestral sites, respectively. Further, ARPIP correctly inferred 83.49% and 87.93% of characters including gaps. Finally, over 99.95% of gap characters were inferred correctly for the two data sets (see Table 2).

Coronavirus Data

The ongoing SARS-CoV-2 pandemic strongly affects our lives, causing an immense interest for phylogenetic analyses of the relevant viral molecular sequences. Like in other coronaviruses, the spike protein in SARS-CoV-2 is important for viral entry into host cells. It is also one of the major determining factors of host range (Belouzard et al. 2012; Zhou and Zhao 2020). We therefore used this protein as an example demonstrating ancestral sequence inference.

SARS-CoV-2 is a member of the *Betacoronavirus* genus which also contains the two other recent human coronavirus strains, namely SARS-CoV and MERS-CoV (Lefkowitz et al. 2018). For our analyses, we selected a small set of available protein sequences from this genus (see Table 3 for the exact sequence list).

Analysis of the coronavirus data set.—The MSAs of the coronavirus sequences were inferred using ProPIP (Maiolo et al. 2018) and PRANK phylogeny-aware webserver (Löytynoja 2014). The total length of the reconstructed MSAs was 2002 and 1929 AAs respectively. The phylogenetic trees were reconstructed by ML in PhyML 3.0 (Guindon et al. 2010), using smart model selection on amino acids and SPR tree moves (Lefort et al. 2017) (see Fig. 4). Then, given an MSA and tree we inferred ancestral sequences with ARPIP. The estimated deletion rates for ProPIP and PRANK's MSAs are respectively $\hat{\mu}=0.242$ and $\hat{\mu}=0.210$. Figure 5 summarizes the resulting ASR by ARPIP comparing to FastML on MSA produced by ProPIP while the results for PRANK can be found in the Appendix S2 of the Supplementary material available on Dryad.

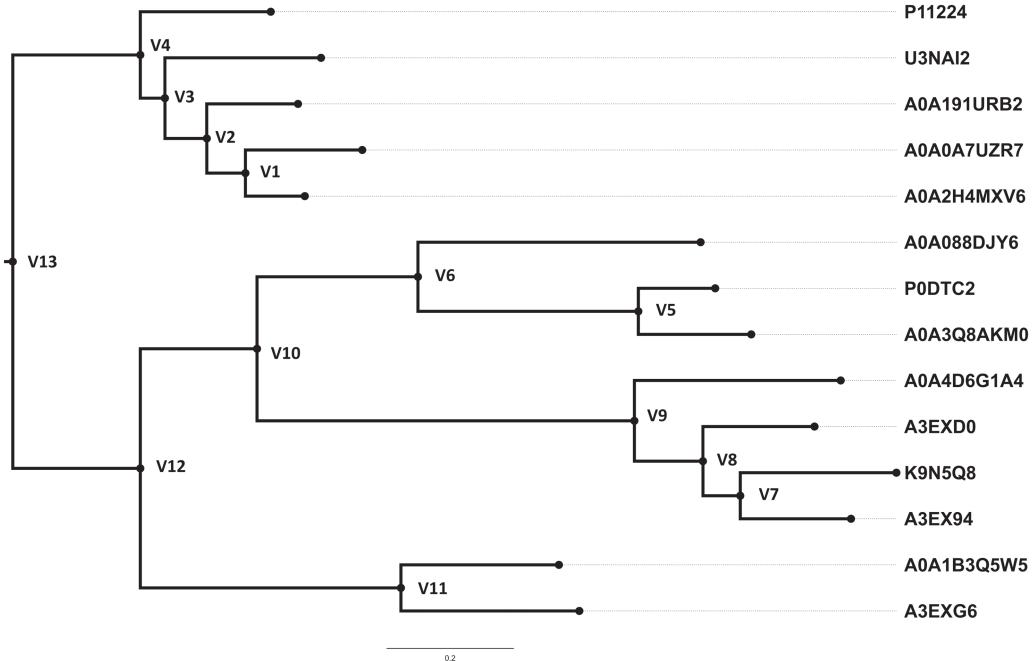


FIGURE 4. Illustration of the rooted *Betacoronavirus* phylogenetic tree which was reconstructed by PhyML 3.0 from the ProPIP alignment. Note that the original tree was unrooted which ARPIP used midpoint rooting method to make the tree rooted.

Comparison Against the State-of-the-Art Methods

At this moment, the two most frequently used ML joint ASR approaches are PAML (Yang 1997) and FastML (Pupko et al. 2000; Ashkenazy et al. 2012), both of which work in linear time with respect to the number of sequences. An important distinction among the methods lies in the way they handle gaps in the alignment. PAML can either ignore gaps in the alignment by removing all columns containing at least one gap character (option used here), or treat all gap characters as ambiguous. For this study, we used PAML on an MSA without any gap characters to compare the accuracy of ancestral character reconstruction. FastML web service (Ashkenazy et al. 2012) uses an ad hoc indel coding (Simmons and Ochoterena 2000) approach to account for indels spanning multiple adjacent characters. Indel coding is done as a separate step in the inference process, done independently from the ancestral state reconstruction. We compare the performance of ARPIP and FastML on an MSA with gaps.

On both simulated data sets, the accuracy of ARPIP appears similar to FastML (see Tables 1 and 2 and Figs. 6, 7, and 8). For example, both algorithms inferred the ancestral state accurately in certain regions (e.g., Figs. 6a, 7a,d, and 8a) and falsely in other regions (e.g., Figs. 6b, 7b,c, and 8b). In Figures 6c and 8d, FastML inferred a character state even though there is no ancestral character, since the insertion happened at the leaf. In certain regions FastML could not determine which internal node had the information (e.g., Fig. 6d), while ARPIP was capable of determining the character position accurately. ARPIP outperformed FastML in determining the gap position in certain regions (e.g.,

Figs. 6c,d, 7b, and 8c,d). On CoV data, the situation was similar meaning FastML could not detect the insertion location (see Fig. 5a) while in conserved regions the inferred states were almost identical (see Fig. 5b).

We also considered scenarios without indels, allowing the evolutionary process to work only through substitutions. In this case the alignments have no gaps, that is, no deletion ($\mu = 0$) and all insertions happen at the root of the tree. Under these conditions, all the algorithms perform reasonably as presented in Figure 9.

DISCUSSION AND CONCLUSION

In this article, we present a one-of-a-kind approach for fast likelihood-based ASR with insertions and deletions. Unlike previous approaches, our method relies on an explicit model of indel and character evolution and allows us to infer the full history of sequence evolution, including insertion and deletion points on a phylogeny. The method is implemented in the probabilistic framework and is based on likelihood calculations under the PIP model. Likelihood computations under this model have linear time complexity with respect to the number of sequences, meaning that our method is highly efficient on large data sets.

We show that on PIP simulated data sets, ARPIP correctly infers at least 95% of indel events and at least 88% of ancestral characters. On the INDELible simulated data, ARPIP correctly infers 83% and 87% of ancestral characters including gaps for low and high indel rates, respectively. ARPIP also correctly places gaps in over 99.95% cases, showing the credibility of our IndelPoints

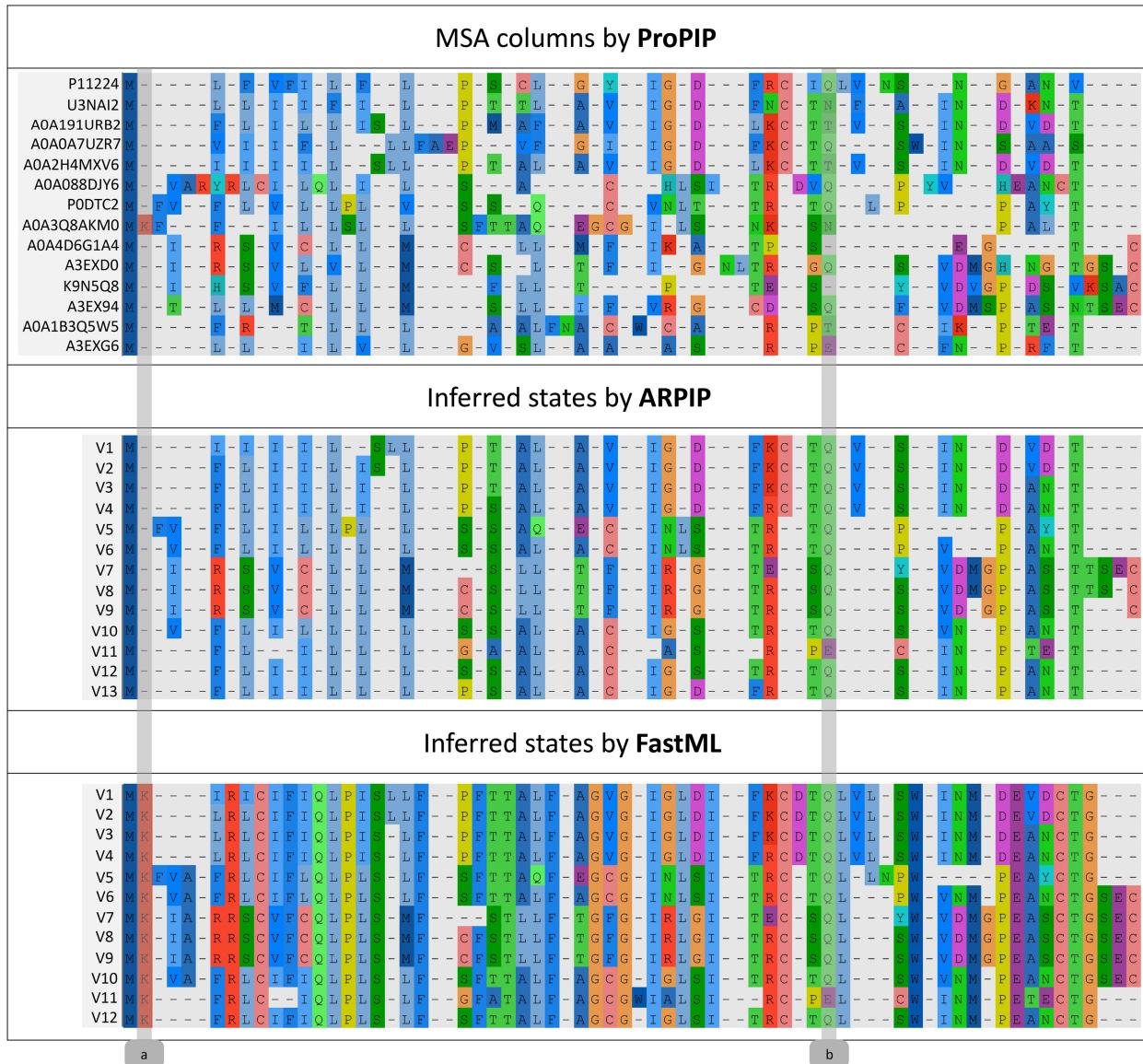


FIGURE 5. Illustration of a snippet from the CoV data set containing the MSA inferred by ProPIP and the ancestral sequences predicted by ARPIP and FastML. a) The region in which ARPIP infers a very different ancestral history, probably due to inferring the insertion point prior to ancestral character inference. FastML inferred no gap in this column perhaps due to the adjacent (first) column. b) The region in which both algorithms had similar inferences of the ancestral states. Note that FastML algorithm works on an unrooted tree which, compared to ARPIP, resulted in one fewer internal sequence reconstructed (due to the absence of the root node).

algorithm. For all data sets, we illustrate the performance of our approach in comparison to FastML on alignments with gaps. In addition, we use gapless alignments to compare ARPIP inferences with those by PAML and FastML, showing that our approach performs just as well for data without gaps.

While indel events represent a major mutational process of gene evolution (Söding and Lupas 2003), they are rarely accounted for in ASR. ARPIP expands the reconstruction possibilities to include more divergent and gappy sequences allowing us to study a wider range of resurrected ancestral molecules, investigating the functional importance of indels in ancestral proteins. This is particularly valuable for proteins separated by

large divergences or within “more flexible” loop regions, as indels frequently occur in regions where amino acid sequences are not well conserved (Taylor et al. 2004a).

While single residue indel modeling may be viewed as a limitation, certain types of genetic material exhibit specifically these kind of indel events more often than others. For example, single nucleotide indels are predominant between recently diverged DNA sequences from various organisms (Tao et al. 2007) and in non-coding DNA sequences (Yamane et al. 2006). While most ASR is done on coding sequences to investigate the properties of reconstructed proteins, it has recently been shown that many trait-associated loci, including some associated with disease, lie outside protein-coding

MSA with gappy columns																	
L1	-	-	Q	G	-	F	-	D	E	-	-	Q	E	L	W		
L2	-	-	E	E	K	L	-	D	E	L	-	-	L	Q	T	Y	
L3	-	-	G	T	S	-	-	K	-	I	V	-	L	S	I	K	
L4	-	-	E	E	-	W	-	D	E	-	-	A	L	V	F	R	
L5	-	-	D	T	-	W	-	D	T	-	-	A	V	L	L	N	
L6	-	-	E	E	-	W	-	D	Q	-	-	S	P	V	F	R	
L7	-	-	E	D	-	K	N	Q	-	P	-	P	N	M	A	N	
L8	T	-	D	-	F	D	Q	-	V	-	E	K	-	T	V	F	A
True simulated states																	
V1	-	-	Q	E	-	F	-	D	E	-	-	Q	E	M	W		
V2	-	-	E	E	-	W	-	D	Q	-	-	V	-	L	V	F	
V3	-	-	E	E	-	W	-	D	Q	-	-	D	-	I	L	A	
V4	-	-	E	E	-	W	-	D	Q	-	-	V	-	L	T	I	
V5	T	-	E	E	-	W	-	D	Q	-	-	P	V	F	A	N	
V6	-	-	E	E	-	W	-	D	Q	-	-	V	-	L	T	I	
V7	-	-	E	E	-	W	-	D	Q	-	-	V	-	L	T	I	
Inferred states by ARPIP																	
V1	-	-	Q	G	-	F	-	D	E	-	-	Q	E	L	W		
V2	-	-	E	E	-	W	-	D	Q	-	-	A	L	V	F		
V3	-	-	E	E	-	W	-	D	Q	-	-	D	-	I	L	A	
V4	-	-	E	E	-	W	-	D	Q	-	-	V	-	L	H	I	
V5	T	-	D	-	W	D	Q	-	V	-	E	K	-	P	V	F	
V6	-	-	E	E	-	W	-	D	Q	-	-	V	-	L	H	I	
V7	-	-	E	E	-	W	-	D	Q	-	-	V	-	L	H	I	
Inferred states by FastML																	
V1	T	-	Q	G	K	F	K	D	E	L	-	-	V	I	H	M	
V2	T	-	E	E	K	W	K	D	E	L	-	-	D	H	I	L	
V3	T	-	E	E	K	W	K	D	E	L	-	-	D	H	I	L	
V4	T	-	E	E	K	W	K	D	E	L	-	-	V	I	H	M	
V5	T	G	D	E	-	W	K	D	Q	L	-	-	P	V	F	A	
V6	T	-	E	E	K	W	K	D	Q	L	-	-	V	I	H	M	

FIGURE 6. A snippet from the PIP simulated data set containing the true simulated MSA and ancestors and the ancestral sequences predicted by ARPIP and FastML. a) A region where both ARPIP and FastML accurately inferred the ancestral states. b) A region where both algorithms estimated the ancestral character incorrectly. c) A region where FastML inferred ancestral characters even though there were none in the simulation. d) A region where there was a single ancestral character but FastML inferred its position incorrectly. Note that FastML algorithm works on an unrooted tree which, compared to ARPIP, resulted in one fewer internal sequence reconstructed (due to the absence of the root node).

MSA with gappy columns																		
L1	A	I	R	G	A	A	A	Q	T	R	R	P	G	A	G	E	G	K
L2	C	T	W	C	A	A	O	P	Q	G	P	A	G	G	T	G	R	
L3	M	V	T	G	T	L	D	V	K	N	K	P	K	T	G	R		
L4	I	V	T	G	A	Y	A	K	C	R	G	P	K	R	G	D	G	
L5	T	V	H	G	A	Y	A	V	I	R	V	V	K	A	S	P	R	
L6	T	V	K	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	
L7	T	E	L	E	A	M	N	R	H	R	H	P	G	D	A	Y	G	
L8	V	I	R	G	N	M	D	D	E	R	D	E	O	G	E	-	G	
True simulated states																		
V1	A	I	R	G	A	A	A	Q	T	R	R	P	G	A	G	E	G	K
V2	I	V	T	G	A	Y	A	K	C	R	P	K	R	G	D	G	G	
V3	T	V	S	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	
V4	T	V	S	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	
V5	T	L	S	G	A	Y	G	P	H	K	D	E	G	A	Y	G	G	
V6	T	V	S	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	
V7	T	V	S	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	
Inferred states by ARPIP																		
V1	A	I	R	G	A	A	A	Q	T	R	R	P	G	A	G	E	G	K
V2	I	V	T	G	A	Y	A	K	C	R	P	K	R	G	D	G	G	
V3	T	V	R	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	
V4	T	V	R	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	
V5	T	L	R	G	A	Y	G	P	H	K	D	E	G	A	Y	G	G	
V6	T	V	R	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	
V7	T	V	R	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	
Inferred states by FastML																		
V1	A	I	R	G	A	A	A	Q	T	R	R	P	G	A	G	E	G	K
V2	I	V	T	G	A	Y	A	K	C	R	P	K	R	G	D	G	G	
V3	T	V	R	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	
V4	T	V	R	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	
V5	T	L	R	G	A	Y	G	P	H	K	D	E	G	A	Y	G	G	
V6	T	V	R	G	A	Y	A	T	R	K	P	N	R	G	E	A	G	

FIGURE 7. A snippet from the INDELible simulated data set with indel rate 0.01 containing the true simulated MSA, ancestors and the ancestral sequences predicted by ARPIP and FastML. a) A region where both ARPIP and FastML accurately inferred the ancestral states. b) A region that the FastML inferred the gaps incomplete while ARPIP missed the the character state. c) A region where both algorithms estimated

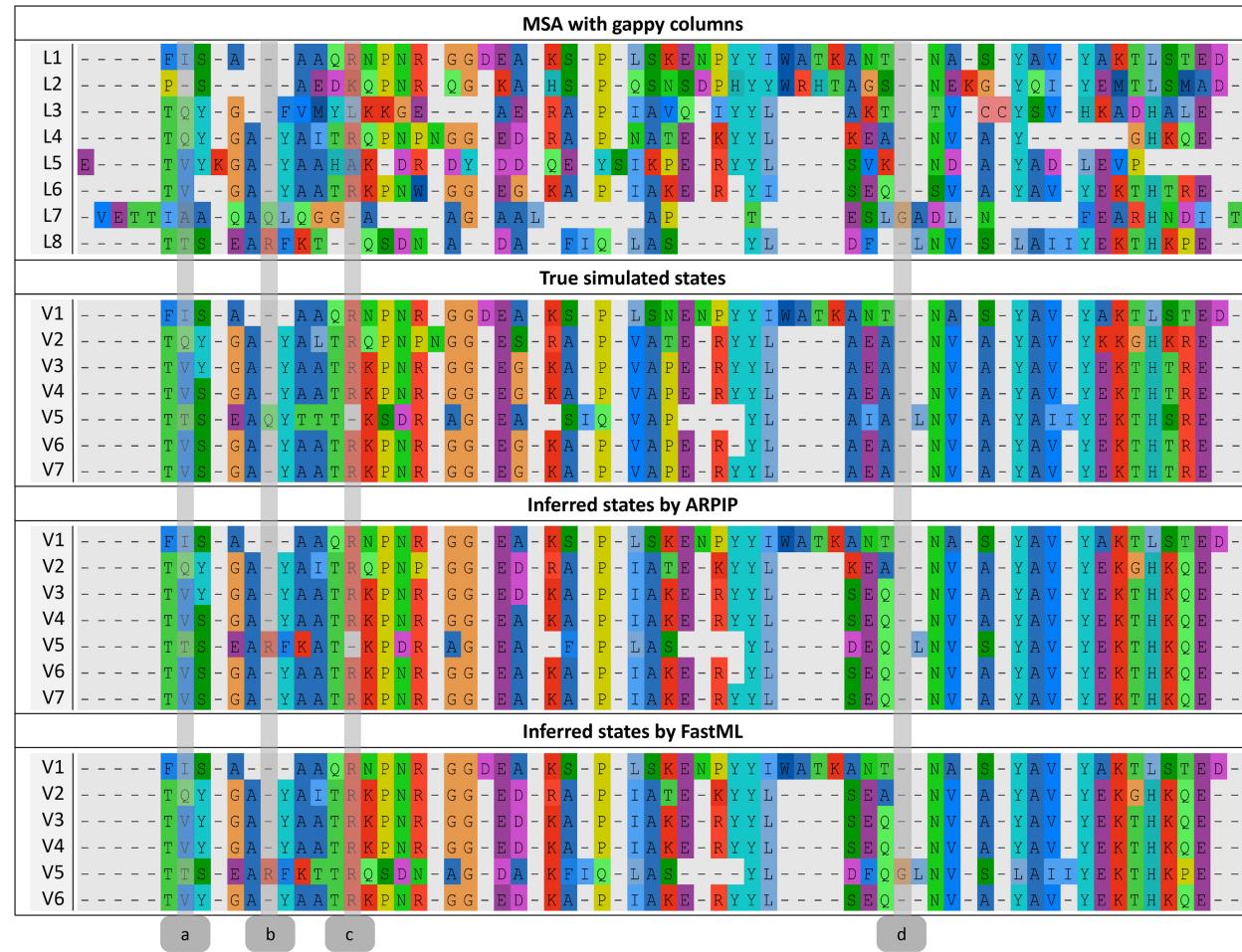


FIGURE 8. A snippet from the INDELible simulated data set with indel rate 0.05 containing the true simulated MSA and ancestors and the ancestral sequences predicted by ARPIP and FastML. a) A region where both ARPIP and FastML accurately inferred the ancestral states. b) A region where both algorithms estimated the indel events correctly but the ancestral character incorrectly. c) A region where FastML missed the gap character but ARPIP inferred it correctly. d) A region where FastML inferred ancestral characters even though there were none in the simulation. Note that FastML algorithm works on an unrooted tree which, compared to ARPIP, resulted in one fewer internal sequence reconstructed (due to the absence of the root node).

regions (Kellis et al. 2014). ARPIP can be used to reconstruct noncoding sequences with meaningful biological assumptions, which could be an additional avenue of exploration for disease-related ASR.

ARPIP paves the way for even more new types of indel analyses. The approach can be expanded to analyze the patterns of insertions and deletions by including rate heterogeneity, for example, allowing us to detect lineage-specific patterns through time. We can include site-specific indel rate variation, allowing us to see the difference in indel evolution in different functional regions of proteins such as loops or active sites. Then, we can investigate the occurrence and consequences of indels in specific regions such as indel-tolerant regions of the genome and relation between gene function and indel frequency (Taylor et al. 2004b). Moreover, in the long run our method can be used to extend and potentially improve more sophisticated probabilistic approaches such as (Groussin et al. 2014), which accounts not only for gene-trees but also for

species history, therefore including gene gain/loss and horizontal transfer in the inference.

While some other methods have attempted to reconstruct exact indel histories, the only other currently existing method in the frequentist framework can only handle small data sets (Diallo et al. 2007). Even though PIP makes simplifying assumptions like site independence, which only allows us to model single residue indels, the explicit evolutionary model makes indel events interpretable. Moreover, as the method has linear time complexity, we can use this approach as a building block in integrated alignment-tree-ancestor inference (Pečerska et al. 2021), using the indel points under PIP as a starting point for integrating more complex models of indel evolution, for example, moving on to long indel models.

Like other likelihood-based approaches, our method in theory allows us to explore both optimal and suboptimal reconstructions in follow-up analyses. It has been argued that a single reconstruction (i.e., a point

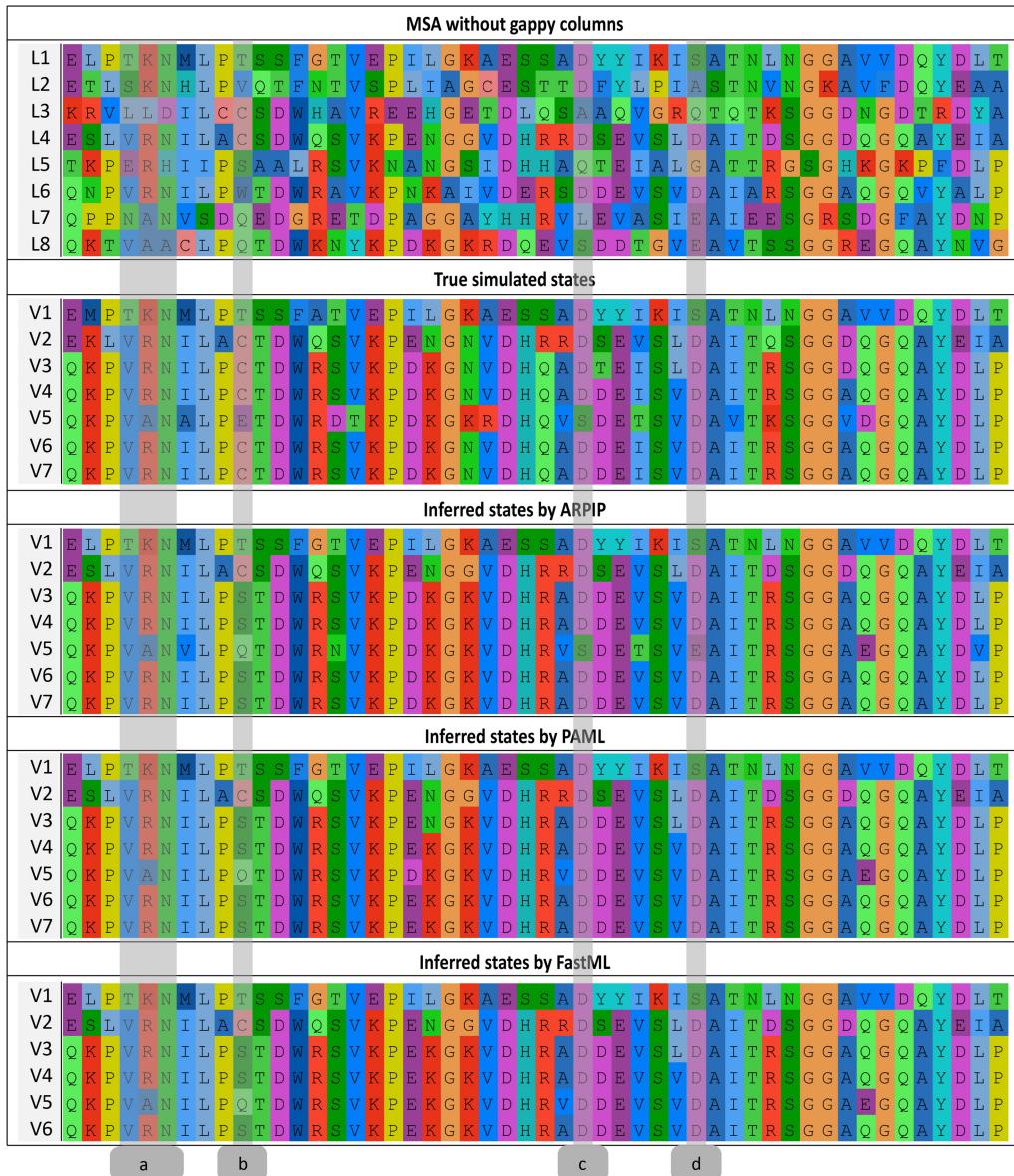


FIGURE 9. A gapless snippet from the PIP simulated data set containing the true simulated MSA and ancestors and the ancestral sequences predicted by ARPIP, PAML and FastML. a) A region where all algorithms accurately inferred the ancestral state. b) A region where all algorithms made mistakes. c) A region where FastML and PAML made incorrect inferences but ARPIP inferred the ancestral state correctly. d) A region where all algorithms accurately inferred the ancestral states except ARPIP. Note that FastML algorithm works on an unrooted tree which, compared to ARPIP, resulted in one fewer internal sequence reconstructed (due to the absence of the root node).

estimate) can be inadequate in cases when the likelihood surface is nonconvex and contains multiple local optima (Joy et al. 2016), which can lead to systematic bias (Yang 2014, p.131). Since ARPIP is in essence an empirical Bayes method, we can extend the method to account for the uncertainty in our estimates by working with probability profiles of characters and gaps rather than inferences fixed to the optimal estimates (Williams et al. 2006).

SUPPLEMENTARY MATERIAL

Data available from the Dryad Digital Repository: <http://dx.doi.org/10.5061/dryad.wstqjq2nj>.

AVAILABILITY OF CODE AND EXPERIMENTAL DATA

The proposed algorithm has been implemented based on Bio++ open source library (Guéguen et al.

2013) using C++ programming language. Our code with a brief user manual is freely available at <https://github.com/acg-team/bpp-ARPIP> under GNU GPLv3 licence. The data underlying this article are also available from the Dryad Digital Repository at <https://doi.org/10.5061/dryad.wstqjq2nj>.

FUNDING

This work was supported by the Swiss National Science Foundation (SNSF) [31003A_176316 to M.A]. The funding body did not play any role in the design of the study and collection, analysis, and interpretation of data and in writing the manuscript.

ACKNOWLEDGMENTS

We would like to thank A. Bouchard-Côté (University of British Columbia) for providing his code JavaPIP to simulate sequences under the PIP, and our master student J. Peechatt for his preliminary work which helped to develop this method.

REFERENCES

- Ashkenazy H., Penn O., Doron-Faigenboim A., Cohen O., Cannarozzi G., Zomer O., Pupko T. 2012. Fastml: a web server for probabilistic reconstruction of ancestral sequences. *Nucleic Acids Res.* 40:W580–W584.
- Bateman A., Martin M.-J., Orchard S., Magrane M., Agivetova R., Ahmad S., Alpi E., Bowler-Barnett E.H., Britto R., Bursteinas B. et al. 2020. Uniprot: the universal protein knowledgebase in 2021. *Nucleic Acids Res.* 49:D480aL–D489.
- Belouzard S., Millet J.K., Licitra B.N., Whittaker G.R. 2012. Mechanisms of coronavirus cell entry mediated by the viral spike protein. *Viruses* 4:1011–1033.
- Bouchard-Côté A. 2010. Probabilistic models of evolution and language change [Ph.D. Thesis]. University of California at Berkeley.
- Bouchard-Côté A., Jordan M.I. 2013. Evolutionary inference via the Poisson indel process. *Proc. Natl. Acad. Sci. USA* 110:1160–1166.
- Brent R.P. 1973. Algorithms for minimization without derivatives. Englewood Cliffs, NJ: Prentice Hall. p. 195.
- Brintnell E., Gupta M., Anderson D.W. 2021. Phylogenetic and ancestral sequence reconstruction of SARS-CoV-2 reveals latent capacity to bind human ACE2 receptor. *J. Mol. Evol.* 89:656–664.
- Chang B.S., Ugalde J.A., Matz M.V. 2005. Applications of ancestral protein reconstruction in understanding protein function: GFP-like proteins. In: Zimmer E.A., Roalson E.H., editors. *Methods in enzymology*, vol. 395. Cambridge (MA): Academic Press. p. 652–670.
- Dessimoz C., Gil M. 2010. Phylogenetic assessment of alignments reveals neglected tree signal in gaps. *Genome Biol.* 11:R37.
- Diallo A.B., Makarenkov V., Blanchette M. 2007. Exact and heuristic algorithms for the indel maximum likelihood problem. *J. Comput. Biol.* 14:446–461.
- Diallo A.B., Makarenkov V., Blanchette M. 2009. Ancestors 1.0: a web server for ancestral sequence reconstruction. *Bioinformatics* 26:130–131.
- Felsenstein J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach. *J. Mol. Evol.* 17:368–376.
- Fletcher W., Yang Z. 2009. Indelible: a flexible simulator of biological sequence evolution. *Mol. Biol. Evol.* 26:1879–1888.
- Groussin M., Hobbs J.K., Szöllősi G.J., Gribaldo S., Arcus V.L., Gouy M. 2014. Toward more accurate ancestral protein genotype-phenotype reconstructions with the use of species tree-aware gene trees. *Mol. Biol. Evol.* 32:13–22.
- Guéguen L., Gaillard S., Boussau B., Gouy M., Groussin M., Rochette N.C., Bigot T., Fournier D., Pouyet F., Cahais V., Bernard A., Scornavacca C., Nabholz B., Haudry A., Dachary L., Galtier N., Belkhir K., Dutheil J.Y. 2013. Bio++: efficient extensible libraries and tools for computational molecular evolution. *Mol. Biol. Evol.* 30:1745–1750.
- Guindon S., Dufayard J.-F., Lefort V., Anisimova M., Hordijk W., Gascuel O. 2010. New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of phylml 3.0. *Syst. Biol.* 59:307–321.
- Joy J.B., Liang R.H., McCloskey R.M., Nguyen T., Poon A.F. 2016. Ancestral reconstruction. *PLoS Comput. Biol.* 12:e1004763.
- Kellis M., Wold B., Snyder M.P., Bernstein B.E., Kundaje A., Marinov G.K., Ward L.D., Birney E., Crawford G.E., Dekker J., et al. 2014. Defining functional DNA elements in the human genome. *Proc. Natl. Acad. Sci. USA* 111:6131–6138.
- Kimura M. 1980. A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences. *J. Mol. Evol.* 16:111–120.
- Lefkowitz E.J., Dempsey D.M., Hendrickson R.C., Orton R.J., Siddell S.G., Smith D.B. 2018. Virus taxonomy: the database of the International Committee on Taxonomy of Viruses (ICTV). *Nucleic Acids Res.* 46:D708–D717.
- Lefort V., Longueville J.-E., Gascuel O. 2017. SMS: smart model selection in PhyML. *Mol. Biol. Evol.* 34:2422–2424.
- Liberles D.A. 2007. *Ancestral sequence reconstruction*. Oxford University Press on Demand.
- Löytynoja A. 2014. Phylogeny-aware alignment with prank. In: Russell D.J., editor. *Multiple sequence alignment methods*. Totowa (NJ): Humana Press. p. 155–170.
- Maiolo M. 2019. Progressive multiple sequence alignment with indel evolution [Ph.D. thesis]. [Lausanne]: University of Lausanne.
- Maiolo M., Zhang X., Gil M., Anisimova M. 2018. Progressive multiple sequence alignment with indel evolution. *BMC Bioinformatics* 19:331.
- Pagel M. 1999. Inferring the historical patterns of biological evolution. *Nature* 401:877–884.
- Pečerska J., Gil M., Anisimova M. 2021. Joint alignment and tree inference. *bioRxiv*. Cold Spring Harbor Laboratory.
- Pupko T., Pe I., Shamir R., Graur D. 2000. A fast algorithm for joint reconstruction of ancestral amino acid sequences. *Mol. Biol. Evol.* 17:890–896.
- Simmons M.P., Ochoterena H. 2000. Gaps as characters in sequence-based phylogenetic analyses. *Syst. Biol.* 49:369–381.
- Söding J., Lupas A.N. 2003. More than the sum of their parts: on the evolution of proteins from peptides. *Bioessays* 25: 837–846.
- Starr T.N., Zepeda S.K., Walls A.C., Greaney A.J., Alkhovsky S., Veesler D., Bloom J.D. 2022. Ace2 binding is an ancestral and evolvable trait of sarbecoviruses. *Nature* 603:913–918.
- Tao S., Fan Y., Wang W., Ma G., Liang L., Shi Q. 2007. Patterns of insertion and deletion in mammalian genomes. *Curr. Genomics* 8:370–378.
- Taylor M.S., Ponting C.P., Copley R.R. 2004a. Occurrence and consequences of coding sequence insertions and deletions in mammalian genomes. *Genome Res.* 14:555–566.
- Taylor M.S., Ponting C.P., Copley R.R. 2004b. Occurrence and consequences of coding sequence insertions and deletions in mammalian genomes. *Genome Res.* 14:555–566.
- Thorne J.L., Kishino H., Felsenstein J. 1991. An evolutionary model for maximum likelihood alignment of DNA sequences. *J. Mol. Evol.* 33:114–124.
- Thorne J.L., Kishino H., Felsenstein J. 1992. Inching toward reality: an improved likelihood model of sequence evolution. *J. Mol. Evol.* 34:3–16.
- Thornton J.W. 2004. Resurrecting ancient genes: experimental analysis of extinct molecules. *Nat. Rev. Genetics* 5:366–375.
- Whelan S., Goldman N. 2001. A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach. *Mol. Biol. Evol.* 18:691–699.
- Williams P.D., Pollock D.D., Blackburn B.P., Goldstein R.A. 2006. Assessing the accuracy of ancestral protein reconstruction methods. *PLoS Comput. Biol.* 2:e69.
- Yamane K., Yano K., Kawahara T. 2006. Pattern and rate of indel evolution inferred from whole chloroplast intergenic regions in sugarcane, maize and rice. *DNA Res.* 13:197–204.
- Yang Z. 1997. Paml: a program package for phylogenetic analysis by maximum likelihood. *Bioinformatics* 13:555–556.
- Yang Z. 2007. Paml 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* 24:1586–1591.
- Yang Z. 2014. *Molecular evolution: a statistical approach*. Oxford: Oxford University Press.
- Yang Z., Kumar S., Nei M. 1995. A new method of inference of ancestral nucleotide and amino acid sequences. *Genetics* 141:1641–1650.
- Zakas P.M., Brown H.C., Knight K., Meeks S.L., Spencer H.T., Gaucher E.A., Doering C.B. 2017. Enhancing the pharmaceutical properties of protein drugs by ancestral sequence reconstruction. *Nat. Biotechnol.* 35:35.
- Zhou G., Zhao Q. 2020. Perspectives on therapeutic neutralizing antibodies against the novel coronavirus SARS-CoV-2. *Int. J. Biol. Sci.* 16:1718.