

Integrated *de novo* gene prediction and peptide assembly of metagenomic sequencing data

Sirisha Thippabhotla¹, Ben Liu¹, Adam Podgorny², Shibu Yooseph³, Youngik Yang⁴, Jun Zhang^{5,6} and Cuncong Zhong^{1,*}

¹Department of Electrical Engineering and Computer Science, The University of Kansas, Lawrence, KS 66045, USA, ²Center for Computational Biology, The University of Kansas, Lawrence, KS 66045, USA, ³Department of Computer Science, Genomics and Bioinformatics Cluster, University of Central Florida, Orlando, FL 32816, USA, ⁴National Marine Biodiversity Institute of Korea, 101-75, Jangsan-ro, Janghang-eup, Seochun-gun, Chungchungnam-do, 33662, South Korea, ⁵Division of Medical Oncology, Department of Internal Medicine, University of Kansas Medical Center, Kansas City, KS 66160, USA and ⁶Department of Cancer Biology, University of Kansas Cancer Center, Kansas City, KS 66160, USA

Received September 22, 2021; Revised December 03, 2022; Editorial Decision February 14, 2023; Accepted February 18, 2023

ABSTRACT

Metagenomics is the study of all genomic content contained in given microbial communities. Metagenomic functional analysis aims to quantify protein families and reconstruct metabolic pathways from the metagenome. It plays a central role in understanding the interaction between the microbial community and its host or environment. *De novo* functional analysis, which allows the discovery of novel protein families, remains challenging for high-complexity communities. There are currently three main approaches for recovering novel genes or proteins: *de novo* nucleotide assembly, gene calling and peptide assembly. Unfortunately, their information dependency has been overlooked, and each has been formulated as an independent problem. In this work, we develop a sophisticated workflow called integrated Metagenomic Protein Predictor (iMPP), which leverages the information dependencies for better *de novo* functional analysis. iMPP contains three novel modules: a hybrid assembly graph generation module, a graph-based gene calling module, and a peptide assembly-based refinement module. iMPP significantly improved the existing gene calling sensitivity on unassembled metagenomic reads, achieving a 92–97% recall rate at a high precision level (>85%). iMPP further allowed for more sensitive and accurate peptide assembly, recovering more reference proteins and delivering more hypothetical protein sequences. The high performance of iMPP can provide a more comprehensive and unbiased view of the microbial communi-

ties under investigation. iMPP is freely available from <https://github.com/Sirisha-t/iMPP>.

INTRODUCTION

Microbial communities are ubiquitously present in many environmental niches on earth, including soil (1), water (2) and air (3). They are a critical component of the human system, playing important roles in maintaining human health and wellbeing (4–6). Human microbiome dysbiosis can lead to various diseases, such as obesity (7–10), diabetes (11,12) and inflammatory bowel disease (13–15). On the other hand, human microbiome intervention has recently been explored as a meaningful non-invasive treatment. For example, *Salmonella*, *Escherichia*, and *Clostridium* are used as anticancer agents with highly promising effects in cancer therapeutics (16–18). Certain microbes are correlated with the response and toxicity from cancer treatments (19,20). Advances in next-generation sequencing (NGS) enable metagenomics, the study of the genomic content of a microbial community as a whole (21,22). Metagenomic sequencing data allows one to examine the taxonomic composition of the microbial community (23–25). More importantly, it further enables protein family profiling (26–28) and metabolic pathway reconstruction (29,30). This information is critical to unlocking the functional potential of the microbial community and elucidating its interactions with the environment.

Metagenomic functional analysis usually begins with homology search, such as aligning the sequencing reads against functionally annotated genomes (e.g. NCBI RefSeq) or protein databases (NCBI NR or UniProt (31)) using BLAST (32). However, due to the incompleteness of current databases, this approach may overlook functional elements encoded by previously-unseen microbial species and novel

*To whom correspondence should be addressed. Tel: +1 785 864 0739; Email: cczhong@ku.edu

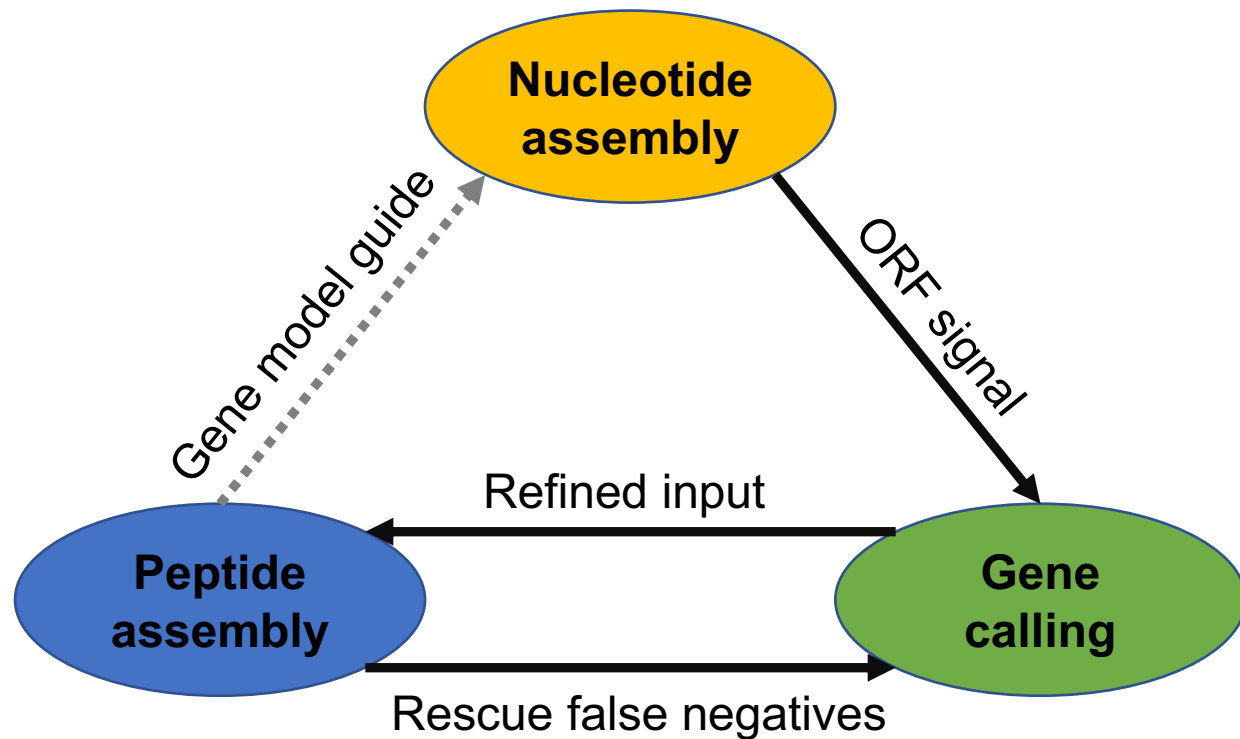


Figure 1. The informational dependency among *de novo* nucleotide assembly, gene calling and peptide assembly. Solid arrows indicate dependencies that have been utilized by iMPP to improve metagenomic functional annotation.

protein families, yielding a biased view of the community's function. Alternatively, a reference-independent approach first assembles the sequencing reads into complete or near-complete genome sequences using *de novo* genome assemblers such as Meta-IDBA (33), MEGAHIT (34), MetaVelvet (35,36) and metaSPAdes (37). Then, it attempts to find open reading frames (ORFs) directly from the assembled genomes based on signals such as gene length, GC-content, and codon usage that characterize most protein-coding genes. The so-called *de novo* gene calling step can be handled by software packages like Glimmer (38), GeneMark (39), and Prodigal (40). When long enough genomic sequences with stable and complete ORF signals are available, *de novo* gene calling is often reliable. However, this step becomes more challenging on fragmented sequences (e.g. unassembled reads). More sophisticated computational models and algorithms are thus required to solve this problem. Software packages that support fragmented gene calling include MetaGeneAnnotator (41), FragGeneScan (42), Orphelia (43), Glimmer-MG (44), MetaGeneMark (45) and MetaProdigal (46). Despite being less accurate than their genome-scale counterparts (42,43), fragmented gene callers can detect low-abundance protein-coding reads that are difficult to assemble. They output the detected protein-coding reads, whose corresponding peptide sequences can be further assembled into peptide contigs using *de novo* peptide assemblers such as SPA (47,48), PLASS (49) and MetaPA (50).

The three *de novo* functional analysis approaches discussed above, i.e. *de novo* nucleotide assembly, gene calling, and peptide assembly, strongly depend on each other (Figure 1). First, nucleotide assembly reconstructs longer

genomic sequences with stronger and more stable ORF signals, which is expected to improve gene calling (42,43). For example, Graph2Pro (51) explicitly couples nucleotide assembly and gene calling by searching ORFs from paths in the nucleotide assembly graph. Second, gene calling can benefit downstream peptide assembly by providing refined short peptide sequences as input. The peptide assembler SPA (47) showed a higher performance when fed with peptide sequences predicted by FragGeneScan (42) compared to those predicted by MetaGeneAnnotator (41). We will further show (in this work) that peptide assemblers that accept all six-frame translations as input can also benefit from a refined input set. Conversely, peptide assembly explores the overlap information among the input short peptide sequences and can improve gene calling by rescuing false-negative predictions. Specifically, if a candidate ORF significantly overlaps with other peptides and is assembled into a long-enough contig, the candidate ORF is likely to be correct. The peptide overlap information is independent of the traditional ORF signals (e.g. codon frequency) and can further contribute to gene calling. Finally, peptide assembly reconstructs longer peptide contigs or even complete protein sequences that can serve as guides to nucleotide assembly. The so-called gene-centric assembly demonstrates better performance than its model-free counterparts (52–54).

Despite the strong informational connection and dependency of *de novo* nucleotide assembly, gene calling, and peptide assembly in metagenomic functional analysis, they have largely been considered and solved independently. Examples include many dedicated metagenome assemblers (33–37), dedicated metagenomic gene callers (41,42,44,45), and dedicated metagenomic peptide assemblers (47–50).

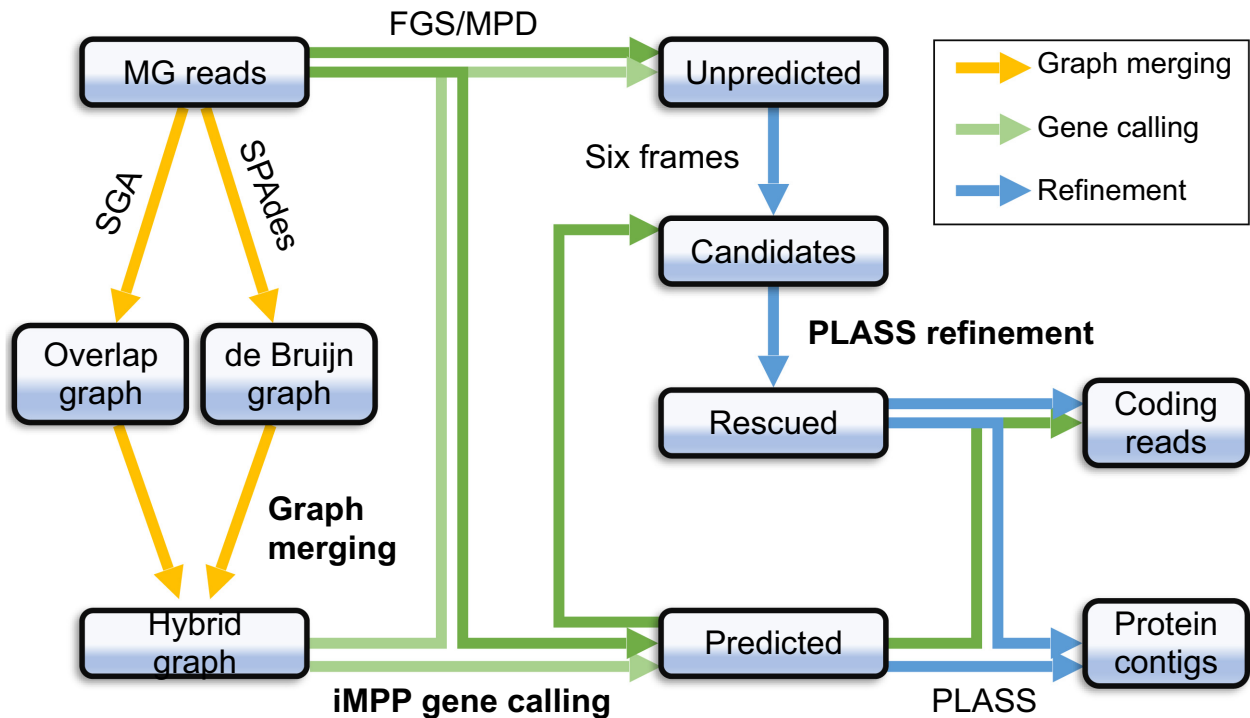


Figure 2. The iMPP workflow overview. Yellow arrows indicate nucleotide assembly information flow; green arrows indicate gene calling information flow using FGS or MPD as the core gene-caller; and blue arrows indicate peptide assembly information flow. The bolded operations, i.e. ‘graph merging’, ‘iMPP gene calling’, and ‘PLASS refinement’ are unique contributions of iMPP and discussed in detail in the Materials and Methods section. ‘MG reads’ stands for metagenomic reads.

While Graph2Pro (51) explicitly couples nucleotide assembly with gene calling, it expects metaproteomic data to validate its protein prediction and lacks a peptide assembly component. To the best of our knowledge, no functional annotation method exists that considers the information dependency among these three approaches and integrates them into a single functional analysis framework. It remains unclear whether doing so is feasible and by how much it can improve metagenomic functional analysis.

We integrate nucleotide assembly, gene calling and peptide assembly into a *de novo* metagenomic functional analysis workflow called integrated Metagenomic Protein Predictor (iMPP). Instead of being a simple sequential execution, iMPP is empowered with three novel modules to fully leverage the information dependency. iMPP constructs a hybrid assembly graph by merging a de Bruijn graph and an overlap graph. The de Bruijn graph information increases graph connectedness, while the overlap graph information retains minor sequence variations. It further contains a novel gene calling module that operates on the merged hybrid graph. The gene calling module is computationally efficient by applying heuristics to eliminate unnecessary graph traversals. Finally, iMPP employs a protein reconstruction module with a two-pass peptide assembly, correcting the gene calling results in the first pass and reconstructing peptide contigs in the second pass. Due to computational efficiency concerns, the current implementation of iMPP does not contain a gene-centric nucleotide assembly module that guides nucleotide assembly with the assembled peptide contigs (Figure 1, the broken gray line).

We benchmarked the performance of iMPP in terms of both *de novo* gene calling and peptide sequence assembly on a mock community, four real metagenomic datasets, and one metatranscriptomic dataset from different environments: human gut, soil, marine, cow rumen and sugarcane rhizosphere. While the performance of the existing gene calling methods is already as high as 80–90%, iMPP further improved it by another ~5%, reaching 85–92% of F-measure. For peptide assembly, we further compared iMPP with two other strategies: one as a sequential integration of gene calling and peptide assembly, and the other as peptide assembly alone. Our evaluations using both real and simulated metagenomic datasets showed that iMPP outperformed both strategies in most assembly statistics, including assembly rate, the number of assembled reads, assembled contig length, N50, reference coverage, and specificity. iMPP successfully recovered ~40–3500 more known protein sequences than the second-best method and reconstructed ~400–434 000 more novel peptide sequences over 60aa. Taken together, iMPP has demonstrated the feasibility and benefit of integrating *de novo* nucleotide assembly, gene calling, and peptide assembly in metagenomic functional analysis.

MATERIALS AND METHODS

The iMPP algorithm

iMPP overview. Figure 2 summarizes the iMPP workflow. iMPP first runs FragGeneScan (42) or MetaProdigal (46) on the unassembled metagenomic (MG) or metatranscriptomic (MT) reads to perform fragmented gene calling. In

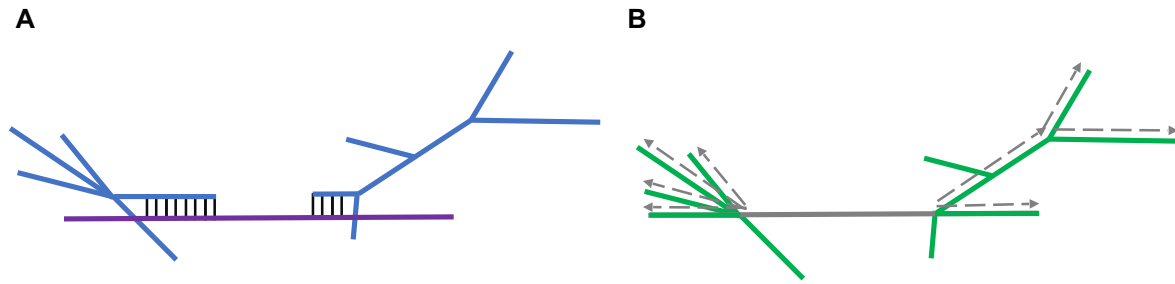


Figure 3. A schematic illustration of the iMPP hybrid graph generation and iMPP gene calling modules. (A) iMPP aligns terminal edges of the assembly overlap graph components (blue) against the de Bruijn graph contigs (purple), and attaches the aligned overlap graph components to the de Bruijn graph contigs. (B) iMPP predicts the protein-coding potential of each edge of the resulted hybrid graph and marks them as either protein-coding (green) or noncoding (gray). iMPP selects the noncoding edges (gray) as anchors and performs depth-first-search from the anchors towards both directions to generate candidate paths.

the rest of this article, we refer to FragGeneScan as FGS and MetaProdigal as MPD for short. In order to leverage sequence overlap information to improve gene calling, iMPP uses nucleotide assemblers SGA (55) and SPAdes (37) to generate assembly overlap graph and de Bruijn graph contigs, respectively. It then merges them into a hybrid graph (see the ‘Assembly Graph Merging’ section). iMPP performs the second pass of gene calling on the edges and paths of the hybrid graph (see the ‘iMPP Gene Calling’ section). Subsequently, iMPP refines the gene calling results by exploiting sequence overlap information among the peptide reads (see the ‘Gene Calling Refinement’ section). Finally, all predicted short peptides are assembled using PLASS (49). Below we focus on the three modules uniquely contributed by iMPP (Figure 2, bolded operations). Note that iMPP is a generic framework that can integrate other options of gene callers, de novo genome assemblers, and de novo peptide assemblers; the options are not limited to FGS, MPD, SGA, SPAdes, and PLASS as have been tested here. More detailed method descriptions, including the chosen parameters and command lines, are available from Supplementary Methods.

Assembly graph merging. iMPP employs a hybrid graph generation module that combines a nucleotide assembly overlap graph and a set of contigs generated by de Bruijn graph assemblers (Figure 3A). de Bruijn graph assembly breaks down the reads into k -mers and models sequence overlap via shared k -mers among reads. It can identify sequence overlaps with a greater sensitivity and often produces more complete assemblies. However, it may overlook minor local sequence variations due to its more aggressive graph simplification strategy. Overlap graph, in contrast, preserves raw sequence variation information but is more fragmentary. Therefore, by merging information from both graphs, we expect to preserve the raw sequence information from the overlap graph and improve the graph connectedness. The idea is similar to hybrid assembly, where longer reads (e.g. PacBio SMRT or Oxford Nanopore MinION) are used to connect short reads (56–58) to improve the overall assembly.

Specifically, iMPP attempts to connect the isolated overlap graph components using de Bruijn graph contigs as a

bridge. iMPP generates an assembly overlap graph using SGA (55) and simplifies the overlap graph by collapsing all unbranched unipaths into single paths (Figure 3A, the blue graph). iMPP then uses SPAdes (under ‘meta’ mode (37)) to generate de Bruijn contigs (Figure 3A, the purple sequence). Denote a vertex with an in- or out-degree of 0 as a *dead end* and an edge containing at least one dead end as a *terminal edge*. iMPP collects all terminal edges from the overlap graph and maps them against all de Bruijn graph contigs. It discards the alignments in which the dead-end sequences are clipped. iMPP then attaches the overlap graph components onto the aligned de Bruijn graph contigs (Figure 3A). iMPP includes all unaligned overlap graph components and de Bruijn graph contigs into the hybrid graph without any modification. This module is similar to the hybrid graph construction module in DRAGoM (59).

iMPP gene calling. Given the hybrid graph, iMPP performs the second pass of gene calling on the paths of the hybrid graph (recall that the first pass of gene calling is performed directly on unassembled reads). Since paths in the hybrid graph contain sequences longer than individual reads, they may contain more complete and stable ORF signals (26,60). However, as the number of paths grows exponentially w.r.t the traversal depth, iMPP employs an ‘anchor and extend’ heuristic to reduce the running time. Specifically, iMPP first runs FGS or MPD on the edges of the hybrid graph. Since microbial genomes are dense in protein-coding genes, the graph usually contains significantly fewer unpredicted edges (i.e. noncoding) than predicted edges. Consequently, iMPP only selects the unpredicted edges as anchors to avoid traversing a large proportion of the graph (Figure 3B). Intuitively, if many predicted edges surround an unpredicted edge, the unpredicted edge is likely to be protein-coding and should also be predicted. iMPP performs a depth-first search (DFS) towards both directions from each anchor (Figure 3B). The DFS terminates after reaching a certain depth, which further bounds the number of paths that need to be reinvestigated. Finally, iMPP re-performs gene calling on the collected paths using FGS (42) or MPD (46). The predicted edges and paths are both considered as protein-coding; the MG reads that can be

Table 1. A summary of benchmark datasets. ‘#Reads’ indicates the total number of reads for the complete dataset; ‘#Genomes’ indicates the number of reference genomes used for subsampling, ‘#Sampled Reads’ corresponds to the number of subsampled reads. ‘Type’ indicates if the sample belongs to Metagenome (MG) or Metatranscriptome data (MT)

Dataset	Accession	Type	Description	Len.	#Reads	#Genomes	#Sampled reads
DS1	SRR341583	MG	Human gut	75	23.4M	3499	9M
DS2	SRR350919	MG	Soil	75	43.9M	65 356	2.9M
DS3	SRR5720229	MG	Marine	150	61.1M	27 216	11.3M
DS4	ERR2027889	MG	Cow rumen	150	109.8M	8980	7.8M
DS5	SRR606249	MG	Mock	102	53.6M	64	49.4M
DS6	SRR14614185	MT	Rhizosphere	101	25.3M	4007	2.3M

mapped to the protein-coding edges and paths are considered as protein-coding reads.

Gene calling refinement. iMPP further refines the gene calling results by utilizing the overlap information in peptide space. Note the difference between this stage and the previous stage, which relies on overlap information in nucleotide space. Due to codon redundancy, reads that cannot be overlapped in nucleotide space (because of synonymous mutations) may be overlapped in peptide space (47). Hence, ORF signals that are missed during nucleotide assembly could be captured by peptide assembly. Specifically, iMPP collects the remaining unpredicted reads and performs all six-frame translations to convert them into *pseudo peptides*. Note that each nucleotide read can associate with up to six pseudo peptides. Then, the pseudo peptides are assembled with the predicted peptides using PLASS (49). Reads with at least one of their pseudo peptides assembled into long-enough contigs are considered as protein-coding; the pseudo peptides contained in the longest contigs are used to determine the frame.

Benchmark datasets

We used six real datasets from different environments (from human gut (61), soil (62), marine (63), cow rumen (64), a mock community, and sugarcane rhizosphere) to benchmark iMPP. We named them DS1–6 respectively. To obtain the ground truth for each dataset, we collected the corresponding taxonomy information from NCBI, and mapped the reads against these reference genomes. The reference genomes and their relative abundances for DS1–6 are available from Supplementary Tables S16–S21, respectively. We compiled all the mapped reads into so-called *subsampling datasets*. We also used the entire set of reads to benchmark the software’s performance on real data; we refer to them as the *complete datasets*. Detailed information is summarized in Table 1 and is available in Supplementary Methods.

We also benchmarked using three simulated datasets, where the first two comprised of reads generated *in silico* from reference genomes, and the third was a CAMI dataset (65). Please see Supplementary Methods and Results as well as Supplementary Table S7 for more information regarding these simulated datasets.

Performance metrics

We benchmarked iMPP with three other strategies in terms of both *de novo* gene calling and peptide assembly. The first

strategy corresponded to the fragmented gene calling directly on the unassembled reads using either FGS or MPD. We refer to the iMPP pipeline that takes FGS as the core gene caller as ‘iMPP(FGS)’, and the MPD-based iMPP as ‘iMPP(MPD)’. The second strategy was to assemble the reads using SGA and then perform gene calling on the assembled contigs. Depending on the gene caller used, we denote this strategy as ‘SGA + FGS’ or ‘SGA + MPD’. The third strategy was similar to the second one, but with SPAdes as the assembler, denoted as ‘SPAdes + FGS’ or ‘SPAdes + MPD’.

We measured the gene calling performance using precision and recall. For the subsampled datasets, we referred to RefSeq (66) for the protein-coding genes in the reference genomes. We defined true positives (TP) as the predicted reads with >60% of their total lengths mapped to the coding regions in the reference genomes, false positives (FP) as the predicted reads that are not mapped to the coding regions, and false negatives (FN) as the unpredicted reads that can be mapped to the coding regions. Then, we computed the recall, precision, and *F*-score as:

$$recall = \frac{TP}{TP + FN}, \quad precision = \frac{TP}{TP + FP},$$

$$F = \frac{2 * recall * precision}{recall + precision}$$

Since no ground truth was available for the complete datasets, we only reported the number of predicted protein-coding reads.

For peptide assembly benchmark, we benchmarked iMPP(FGS) and iMPP(MPD) with two other strategies. The first strategy corresponded to the assembly of FGS (42) or MPD (46) predicted reads using PLASS (49). We refer to this strategy as ‘FGS + PLASS’ or ‘MPD + PLASS’, respectively. This strategy was similar to SPA (47,48), which expected the input to be selected by gene callers. The second strategy was to use the entire set of unfiltered reads, which was the expected input of PLASS. We refer to this strategy as ‘PLASS’. We measured the total number of assembled reads, the number of output contigs, the total length of output contigs, N50, assembly rate (the number of assembled reads over the total number of reads), and chimera rate. To evaluate the correctness of the assembly, we further aligned (using DIAMOND (67)) the contigs against the proteins encoded in the reference genomes (for the subsampled datasets) and the UniProt (31) database (for the complete datasets). A contig was considered true if its aligned proportion was above a certain threshold. We re-

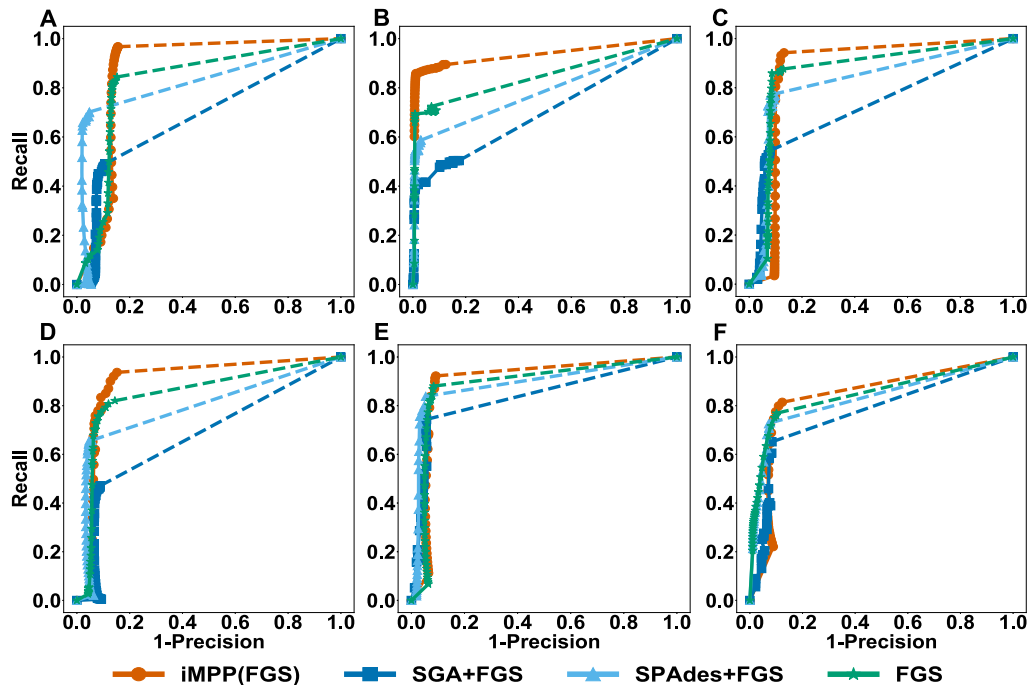


Figure 4. The ROC curves for the de novo gene calling performances of iMPP(FGS), SGA + FGS, SPAdes + FGS and FGS on the six subsampled datasets. (A) DS1, (B) DS2, (C) DS3, (D) DS4, (E) DS5 and (F) DS6.

ported *contig-level specificity* as the total length of the true contigs over the total contig length, and *read-level specificity* as the total number of reads constituting the true contigs over the total number of assembled reads. Finally, to measure sensitivity, we reported reference coverage as the percentage of reference genes covered by the assembled contigs.

All experiments were run on an in-house Linux server equipped with an Intel Xeon E7-4850 CPU and 1T of physical memory. All software with run time reported were executed with 16 threads.

RESULTS

Gene calling benchmark

The gene calling performances of iMPP(FGS) and the other strategies on the six subsampled datasets are summarized in Figure 4. The results were broadly consistent among all datasets, where iMPP(FGS) demonstrated the highest performance, followed by FGS, SPAdes + FGS and SGA + FGS. Specifically, the peak *F*-scores of iMPP(FGS) were between 84.97% and 92.16% among the six datasets (Supplementary Table S1). The second-best strategy, FGS, showed *F*-scores of between 80.03% and 90.24%. iMPP(FGS) showed *F*-score improvement over FGS by 1.92% (DS5)–9.39% (DS2). The minimum improvement was observed from DS5, a low-complexity mock community that contains only 64 microbial species. Given the already-high performance baseline of >80% *F*-score of FGS, the improvement was significant. Strategies that perform gene calling on assembled reads, i.e. SGA + FGS and SPAdes + FGS, performed worse than FGS, potentially be-

cause many reads were not assembled into contigs and were not considered.

Similarly, iMPP(MPD) also showed improvement over MPD on these six datasets (Supplementary Figure S1). iMPP(MPD) showed peak *F*-scores between 83.04% and 92.69%, while MPD was between 77.12% and 90.61% (Supplementary Table S2). iMPP(MPD) outperformed MPD on *F*-score by 2.08% (DS5) and 11.76% (DS2). Overall, the performances of iMPP(FGS) and iMPP(MPD) are highly similar, and they both improved the corresponding core gene callers FGS and MPD. The results shows that the improvement brought by the iMPP framework is likely independent of the core gene caller it integrates.

For the gene calling performance on the complete datasets, we only report the raw prediction counts because no ground truth is available (Figure 5). iMPP(FGS) and FGS predicted more protein-coding reads than the assembly-based strategies SGA + FGS and SPAdes + FGS. This is likely due to the low assembly rate on these datasets. iMPP(FGS) also predicted more reads than FGS, especially on DS1, DS2, and DS6 (14.60%, 18.69% and 15.92% more, respectively). The improvement was marginal on DS3, DS4, and DS5 (2.00%, 3.06% and 5.09% more respectively). The results are consistent with the observations made from the subsampled datasets (Figure 4), where iMPP(FGS) showed the highest recall rate among all strategies. The benchmark results with MPD as the core gene caller also showed a similar trend (Supplementary Figure S2). iMPP(MPD) remained the best method that predicts significantly more protein-coding reads than MPD on DS1, DS2, DS5 and DS6 (18.99%, 11.97%, 13.18% and 32.81% more, respectively) and marginally more on DS3 and DS4 (0.31% and 4.27% more, respectively).

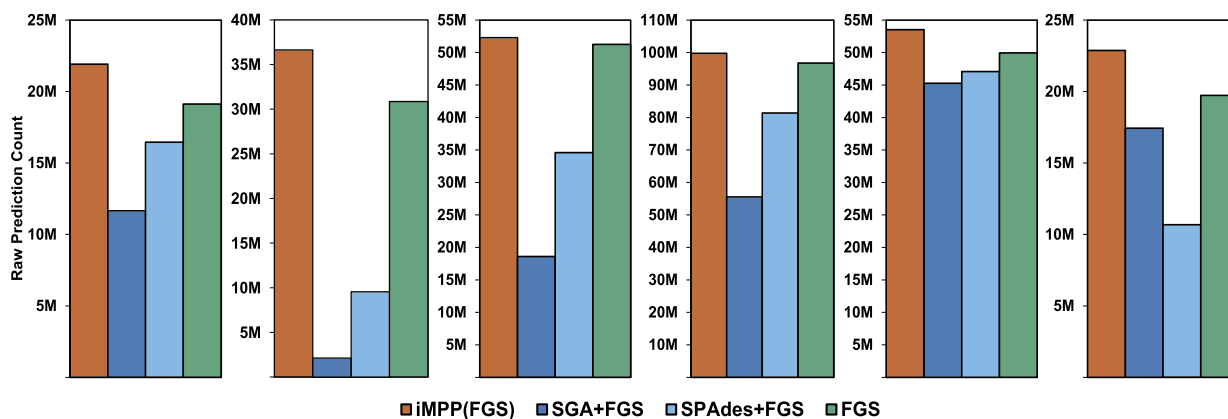


Figure 5. The number of protein-coding reads predicted by iMPP(FGS), SGA + FGS, SPAdes + FGS and FGS on the six complete datasets. Panels from left to right: DS1, DS2, DS3, DS4, DS5 and DS6.

Table 2. Peptide assembly statistics of iMPP(FGS), FGS + PLASS and PLASS on the subsampled datasets. The highest performance in each category is bolded

Dataset	Metrics	iMPP(FGS)	FGS + PLASS	PLASS
DS1 (subsampled)	# Contigs	28 851	25216	27826
	Assembly rate (%)	29.14	29.17	25.93
	Total contig length	2.19M	1.95M	2.09M
	N50(bp)	75	75	75
	Chimera rate (%)	0.033	0.018	0.02
DS2 (subsampled)	# Contigs	13 529	10795	13490
	Assembly rate (%)	8.57	7.16	7.32
	Total contig length	1.21M	956K	1.21M
	N50(bp)	90	89	89
	Chimera rate (%)	0	0	0
DS3 (subsampled)	# Contigs	759K	720K	758K
	Assembly rate (%)	59.97	51.01	58.19
	Total contig length	79.72M	77.08M	79.54M
	N50(bp)	121	118	117
	Chimera rate (%)	0.009	0.012	0.011
DS4 (subsampled)	# Contigs	2.27M	2.23M	2.27M
	Assembly rate (%)	85.54	82.35	79.39
	Total contig length	353.23M	350.73M	353.12M
	N50(bp)	179	177	176
	Chimera rate (%)	0.097	0.099	0.091
DS5 (subsampled)	# Contigs	15.31M	13.95M	15.30M
	Assembly rate (%)	96.20	95.91	95.99
	Total contig length	2.04B	1.88B	2.04B
	N50(bp)	147	142	147
	Chimera rate (%)	2.42	2.23	2.42
DS6 (subsampled)	# Contigs	14 909	13095	14700
	Assembly rate (%)	23.63	18.65	22.7
	Total contig length	1.56M	1.35M	1.55M
	N50(bp)	108	105	107
	Chimera rate (%)	0	0	0

Peptide assembly benchmark

We summarize the peptide assembly benchmark results on the subsampled datasets in Table 2. We only considered peptide contigs that are ≥ 60 aa long. iMPP(FGS) assembled the largest number of contigs and total contig length for all datasets. It outperformed FGS + PLASS by 0.71–27.17% of the total contig length, but with a less significant improvement over PLASS (0.03–4.88%). Note that the peptide assembly module of iMPP only accepted the predicted protein-coding reads as input, which is less than the entire dataset accepted by the PLASS strategy (DS1: 8.65M versus 9.0M, DS2: 2.09M versus 2.9M, DS3: 10.59M ver-

sus 11.3M, DS4: 7.54M versus 7.8M, DS5: 49.29M versus 53.6M and DS6: 2.17M versus 2.3M). However, even with fewer input reads, iMPP(FGS) assembled more contigs in terms of both the quantity and total length. It suggests that eliminating noncoding reads from the input can potentially benefit peptide assembly. On the other hand, iMPP(FGS) also outperformed the FGS + PLASS approach that also refined the input. It suggests that true protein-coding reads should not be excluded from the input or it might harm peptide assembly. By using the most accurate input sets, iMPP(FGS) had the highest assembly rate overall. Although it slightly underperformed FGS + PLASS on DS1

Table 3. The contig- and read-level specificity (%) for the peptide assemblies made by iMPP(FGS), FGS + PLASS, and PLASS at different length thresholds on the subsampled datasets. The highest performance in each category is bolded

Dataset	Len	iMPP(FGS)		FGS + PLASS		PLASS	
		contig	read	contig	read	contig	read
DS1 (subsampled)	60%	83.69	89.90	79.29	86.00	83.60	89.81
	70%	83.39	89.43	79.00	85.61	83.30	89.39
	80%	82.77	88.66	78.36	84.88	82.69	88.55
	90%	81.41	87.08	76.74	83.09	81.36	86.81
DS2 (subsampled)	60%	98.35	99.47	98.44	99.71	98.22	99.43
	70%	98.12	99.21	98.34	99.60	98.02	99.24
	80%	97.87	98.75	97.98	99.11	97.73	98.92
	90%	97.15	98.22	97.05	98.14	96.50	97.95
DS3 (subsampled)	60%	73.12	69.20	71.99	66.31	69.89	61.35
	70%	72.82	68.35	70.98	63.38	68.83	58.16
	80%	71.23	65.73	68.93	58.96	66.84	53.02
	90%	66.21	61.26	60.77	51.04	58.86	50.86
DS4 (subsampled)	60%	84.24	97.12	83.60	96.89	83.00	96.15
	70%	83.88	96.68	83.34	96.43	82.73	95.66
	80%	82.35	95.48	82.78	95.82	82.16	94.92
	90%	81.29	94.66	81.02	94.64	80.58	93.70
DS5 (subsampled)	60%	89.82	94.96	88.95	94.13	89.60	94.94
	70%	86.52	94.12	86.36	93.82	86.31	94.03
	80%	84.48	92.82	84.02	92.49	84.50	92.75
	90%	82.44	91.08	82.18	90.87	82.13	90.98
DS6 (subsampled)	60%	87.83	95.73	86.39	95.92	87.12	96.18
	70%	86.94	94.74	85.16	95.02	86.87	95.21
	80%	85.43	93.97	84.00	94.29	85.42	93.83
	90%	82.43	92.87	80.44	92.13	82.85	92.67

(~0.03%), iMPP(FGS) showed a significantly higher assembly rate than FGS + PLASS on the rest of the datasets by 0.29% (DS5) to 8.96% (DS3). iMPP (FGS) also consistently showed the highest N50, although it remained similar to the other strategies. Finally, iMPP(FGS) showed a slightly higher chimera rate than the other methods by 0.19% in the worst-case scenario (DS5). As expected, the MPD-based benchmark results showed a similar trend (Supplementary Table S3).

We aligned the resulted contigs against the ground-truth reference proteins to investigate the accuracy of the peptide assemblies. The contig- and read-level specificities of iMPP(FGS), FGS + PLASS and PLASS are summarized in Table 3. All three strategies had similar levels of performance, with most of the differences <3%. iMPP(FGS) showed the highest assembly accuracy on DS1, DS3, DS4 and DS5, while FGS + PLASS was the best for DS2. Both methods performed similarly on DS6. For the comparison between iMPP(MPD), MPD + PLASS and PLASS, iMPP(MPD) showed the highest assembly accuracy on DS3–6, while PLASS was the best for DS1 and MPD + PLASS for DS2 (Supplementary Table S4).

We further calculated the proportion of reference protein sequences recovered by the assemblies generated by difference strategies (Figure 6). iMPP(FGS) consistently showed the highest reference coverages at all sequence length thresholds on all six benchmark datasets. The average improvement over the second-best PLASS strategy was 0.96%. The results were in line with the observation that iMPP(FGS) generated more contigs and total contig length than PLASS (Table 2). Similar improvements were also observed for iMPP(MPD) (Supplementary Figure S3). Taken together, iMPP(FGS) and iMPP(MPD) both improved *de*

novo peptide assembly sensitivity and accuracy on the subsampled datasets.

We also performed similar analyses on the complete datasets. The results summarized in Table 4 were largely consistent with what had been observed for the subsampled datasets (Table 2). Specifically, iMPP(FGS) assembled significantly more contigs (2.29–21.24%) and longer total contig length (1.83–24.03%) than FGS + PLASS, and performed similarly as PLASS (0.00–1.43% more assembled contigs, 0.00–0.69% longer total contig length). iMPP(FGS) consistently showed the highest assembly rate and N50 among all datasets, although with marginal improvements. All strategies had the same low chimera rate. The same conclusion can also be made from the MPD-based benchmark results (Supplementary Table S5).

As we did not have the ground truth reference proteins for the complete datasets, we aligned the assembled peptide contigs against the UniProt database (31) to benchmark assembly accuracy. The corresponding contig- and read-level specificities for iMPP(FGS), FGS + PLASS and PLASS are summarized in Table 5. The results were again consistent with the subsampled datasets, with iMPP(FGS) leading in most of the metrics. All accuracies were lower than those for the subsampled datasets, as the complete dataset may contain more novel proteins that cannot be aligned. Interestingly, the second-most accurate strategy appeared to be PLASS for the complete datasets, unlike FGS + PLASS for the subsampled datasets. The reason could be that the ORF model used by FragGeneScan was trained on known protein families and therefore might miss true protein-coding reads from the novel protein families in the complete datasets. The less complete input further led to fragmentary assemblies, where many short contigs failed to be

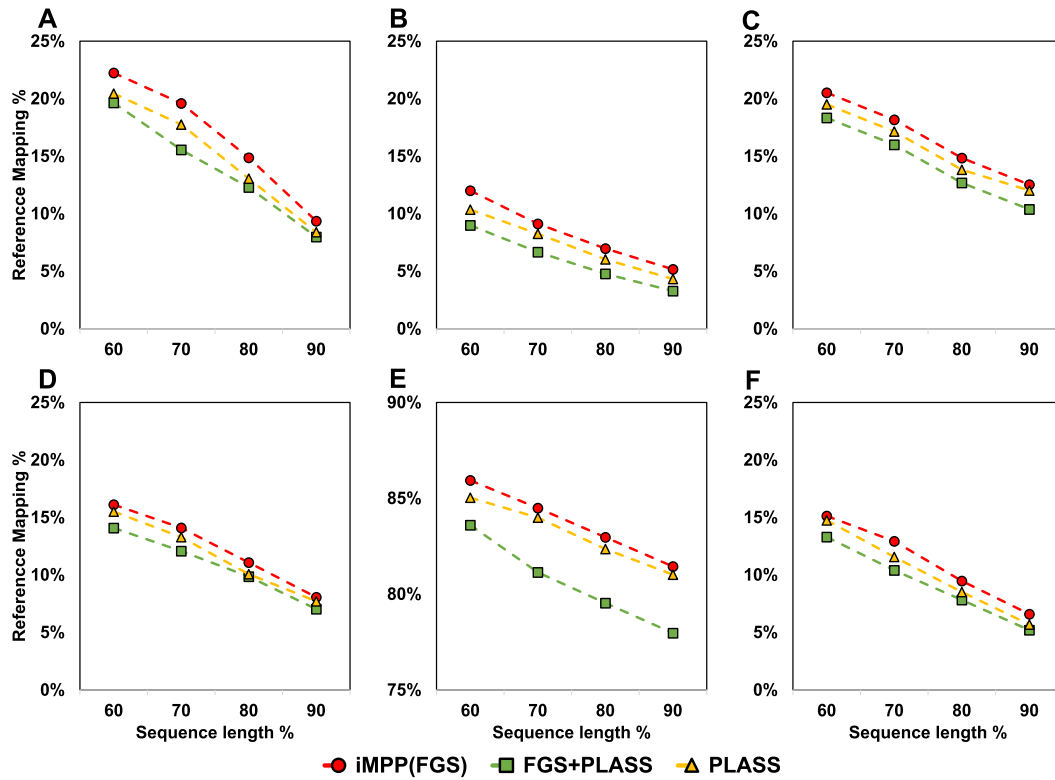


Figure 6. Reference coverages by the peptide contigs assembled by iMPP(FGS), FGS and PLASS on the six subsampled datasets. (A) DS1, (B) DS2, (C) DS3, (D) DS4, (E) DS5 and (F) DS6.

Table 4. Peptide assembly statistics of iMPP(FGS), FGS + PLASS, and PLASS on the complete datasets. The highest performance in each category is bolded

Dataset	Metrics	iMPP(FGS)	FGS + PLASS	PLASS
DS1 (complete)	# Contigs	69 760	66341	69288
	Assembly rate (%)	25.23	25.19	23.56
	Total contig length	5.27M	5.01M	5.23M
	N50(bp)	75	74	74
	Chimera rate (%)	0	0	0
DS2 (complete)	# Contigs	518 842	470011	516296
	Assembly rate (%)	9.89	9.57	8.55
	Total contig length	56.67M	50.81M	56.37M
	N50(bp)	111	110	111
	Chimera rate (%)	0	0	0
DS3 (complete)	# Contigs	7.32M	7.13M	7.32M
	Assembly rate (%)	41.78	39.63	36.27
	Total contig length	884.79M	868.91M	884.79M
	N50(bp)	131	130	131
	Chimera rate (%)	0	0	0
DS4 (complete)	# Contigs	30.86M	30.17M	30.43M
	Assembly rate (%)	66.47	65.79	60.58
	Total contig length	950.88M	880.12M	950.30M
	N50(bp)	205	205	205
	Chimera rate (%)	0	0	0
DS5 (complete)	# Contigs	16.23M	14.83M	16.23M
	Assembly rate (%)	95.91	95.62	95.67
	Total contig length	2.03B	1.99B	2.03B
	N50(bp)	146	141	146
	Chimera rate (%)	2.42	2.23	2.42
DS6 (complete)	# Contigs	172 970	142663	172009
	Assembly rate (%)	32.78	29.49	32.36
	Total contig length	19.39M	15.63M	19.33M
	N50(bp)	115	111	115
	Chimera rate (%)	0	0	0

Table 5. The contig- and read-level specificity (%) for the peptide assemblies made by iMPP(FGS), FGS + PLASS and PLASS at different length thresholds on the complete datasets. The highest performance in each category is bolded

Dataset	Len	iMPP(FGS)		FGS + PLASS		PLASS	
		Contig	Read	Contig	Read	Contig	Read
DS1 (complete)	60%	27.74	50.84	26.19	47.3	26.98	49.17
	70%	27.24	49.99	25.75	46.84	26.21	48.32
	80%	26.24	48.87	24.81	45.97	25.72	47.15
	90%	24.04	48.28	22.02	43.91	25.13	46.77
DS2 (complete)	60%	46.95	54.72	45.78	54.12	45.91	54.4
	70%	46.28	53.95	45.16	53.39	45.22	53.58
	80%	44.83	52.43	44.81	51.94	43.8	52.05
	90%	41.28	49.03	41.15	48.48	40.18	48.68
DS3 (complete)	60%	45.01	54.29	44.82	53.41	44.06	52.58
	70%	44.39	53.71	44.04	52.49	43.29	51.58
	80%	42.96	52.19	42.25	50.65	41.52	49.59
	90%	38.51	49.23	36.57	46.16	35.93	44.78
DS4 (complete)	60%	42.29	42.26	41.72	40.75	41.68	41.29
	70%	41.67	41.13	41.45	40.12	41.42	40.22
	80%	40.88	39.54	40.87	38.24	40.85	38.84
	90%	40.1	36.29	39.42	36.21	39.71	36.02
DS5 (complete)	60%	56.34	62.97	55.10	61.67	55.84	62.44
	70%	53.73	61.22	53.56	59.09	53.93	61.79
	80%	52.31	59.43	52.20	58.91	52.31	59.41
	90%	51.17	57.53	51.12	56.43	51.17	57.38
DS6 (complete)	60%	32.26	31.59	30.8	28.64	26.98	30.1
	70%	31.78	30.65	30.24	28.39	26.46	29.82
	80%	30.44	29.79	29.11	27.74	25.39	29.12
	90%	28.28	28.16	26.44	26.02	23.07	27.26

reliably aligned to references. The results for iMPP(MPD), MPD + PLASS, and PLASS benchmark are available from Supplementary Table S6, where iMPP(MPD) was again the best.

Finally, Figure 7 summarizes the number of peptide contigs that were reconstructed by iMPP(FGS), FGS + PLASS and PLASS and could be aligned to the UniProt database (31), under different reference length thresholds. iMPP(FGS) was able to recover the largest number of known protein sequences from UniProt, followed by PLASS (~40–3500 more peptides). Similarly, iMPP(MPD) predicted ~400–4500 more peptides compared to PLASS (Supplementary Figure S4). These results reconfirmed iMPP's high peptide assembly sensitivity on real datasets.

Benchmark results on simulated datasets

In addition to the real datasets DS1–6, we also benchmarked the iMPP framework on three simulated datasets (Supplementary Table S7). We generated two in-house datasets *in silico*, one from 28 marine microbial genomes and one from 8 *Streptococcus* genomes. We also included the CAMI (65) medium-complexity dataset and subsampled it based on the reference genomes provided by the database. The reference genomes used for subsampling and their relative abundances are summarized in Supplementary Tables S22–S24. More details regarding benchmark results on the simulated datasets can be found from Supplementary Methods and Results, Supplementary Figures S5–S8, and Supplementary Tables S8–S13. For all datasets, because of their relatively low complexity, all methods performed similarly well.

Time-performance tradeoff

We investigated the proportion of true protein-coding reads discovered by different modules of the iMPP(FGS) pipeline (Figure 8 and Supplementary Table S14). Recall that iMPP(FGS) can make ORF predictions in three stages: from the direct application of FGS on unassembled reads, from the iMPP gene calling module on the hybrid graph, and finally from the peptide assembly-based refinement. The most economical way to identify coding reads was to perform fragmented gene calling, as FGS could find >85% of the true positives using ~10% of the total time. The result was consistent with the high performance observed for FGS (42). The remaining ~15% of the protein-coding reads were more challenging to discover, but the majority of them could be discovered using the iMPP gene calling module. It indicates that longer paths from the hybrid assembly graph indeed contain stronger ORF signals and benefit gene calling. However, this module was also the most time-consuming since it performed both overlap graph assembly and de Bruijn graph assembly. It took up ~55–80% of the total runtime of iMPP(FGS). Finally, a very small proportion (2–3%) of the coding reads could be rescued by peptide assembly-based refinement, which took ~15–30% of the total runtime. The same trend is also observed for iMPP(MPD) (Supplementary Figure S9 and Supplementary Table S15).

In comparison between iMPP(FGS) and iMPP(MPD), iMPP(FGS) runs ~10% faster than iMPP(MPD) (Supplementary Tables S14 and S15). However, as shown above, iMPP(FGS) also shows a marginally lower performance. The choice between iMPP(FGS) and iMPP(MPD) can be made based on the available computing resources and the performance need.

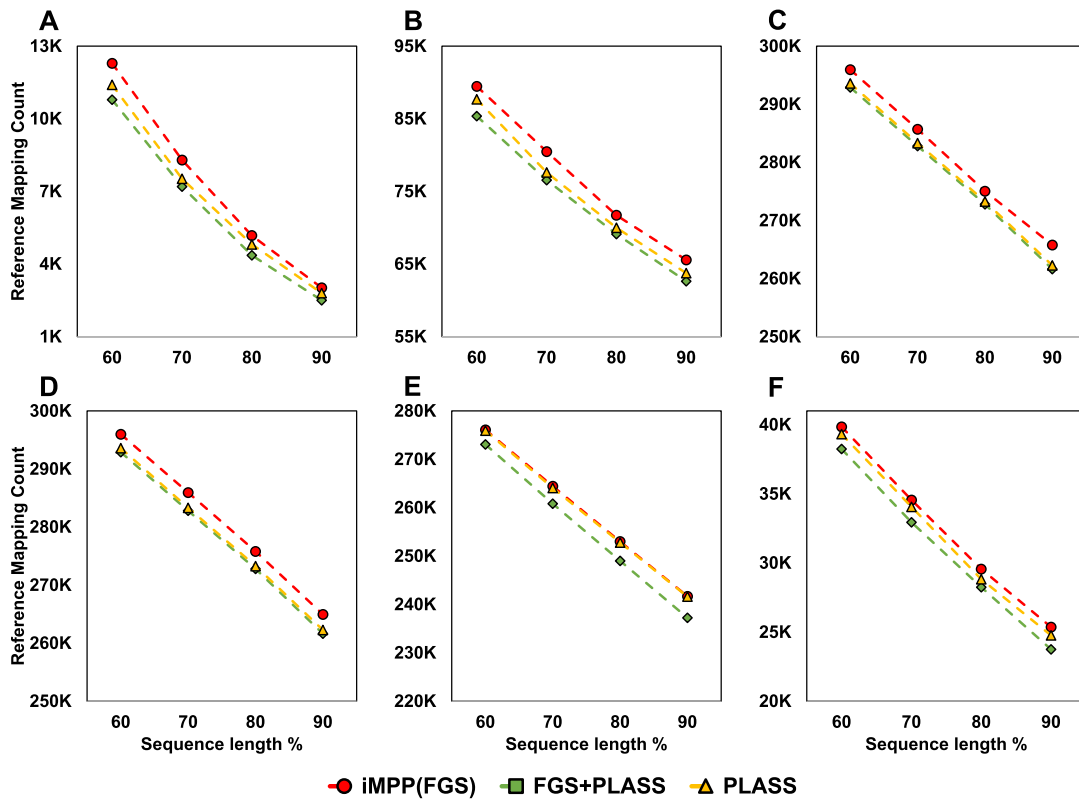


Figure 7. The number of aligned (against UniProt) peptide contigs assembled by iMPP(FGS), FGS and PLASS on the six complete datasets. (A) DS1, (B) DS2, (C) DS3, (D) DS4, (E) DS5 and (F) DS6.

Pathway analysis of the diabetic human gut dataset (DS1)

To show the biological importance of iMPP, we analyzed the diabetic human gut microbiome dataset (DS1, from a Type 2 diabetes patient). We compared the metabolic pathways reconstructed from the genes predicted by iMPP(FGS) and FGS. We used MetaCyc (68) as the pathway database and MinPath (69) as the pathway reconstruction tool (see Supplementary Materials and Methods). MinPath reconstructed 196 more pathways from iMPP(FGS)'s prediction in comparison to FGS's prediction (see Supplementary Table S25). Among these unique pathways, some of them are universal metabolic pathways found in human gut, and some may relate to the changes induced by diabetes.

With iMPP(FGS)'s prediction, MinPath revealed some pathways that are expected to universally exist in human gut environment but were missed by FGS's prediction. For example, the sulfate activation for sulfonation pathway (MetaCyc ID: PWY-5340) and the spermidine biosynthesis pathway (MetaCyc ID: BSUBPOLYAMSYN-PWY). Sulfate activation is a prerequisite step for the degradation of inorganic sulfate by the sulfate-reducing bacteria (SRB), a well-known integral part of the intestinal microbiota (70). Meanwhile, sulfonation integrates a sulfur-containing group to an organic molecule, and the process is used in the biosynthesis of sulfur-containing essential amino acids such as cysteine (71) and sulfated cell-surface glycosaminoglycans that are used as receptors by nearly all kinds of bacteria (72).

The second pathway synthesizes spermidine, a major type of polyamine that is found in both eukaryotic and prokaryotic cells to carry out a wide range of essential biological functions such as gene regulation, stress resistance, and cell proliferation (73). Studies have shown that polyamine found in human lower intestinal tract is primarily synthesized by gut microbiota (74,75). Specifically, spermidine biosynthesis has been confirmed in the human gut bacteria *Bacteroides thetaiotaomicron* and *Fusobacterium varium* (76).

Interestingly, some pathways such as the betaine biosynthesis pathway (MetaCyc ID: PWY-4021) and the ectoine biosynthesis pathway (MetaCyc ID: PWY-4021), were uniquely reconstructed from iMPP(FGS)'s prediction, and appeared to relate to the osmolarity change in the gut environment induced by diabetes. Many betaines serve as organic osmolytes, which are used by cells in response to osmotic stress and to protect them from dehydration (77). While the primary source of betaines is known to come from food, Koistinen *et al.* recently showed that betaine can also be synthesized by the gut microbiota using both mouse and *in vitro* human gut models (78). As a result, the identification of betaine biosynthesis pathway may indicate the prevalence of microbes that can protect themselves from the osmolarity change in the diabetic gut environment, by synthesizing betaine compounds. Similarly, ectoine is also known to be synthesized by the gut microbiota in response to osmolarity stress (79,80). These osmolyte transporter biosynthesis pathways tend to be more active under the diabetic gut environment, believed to buffer

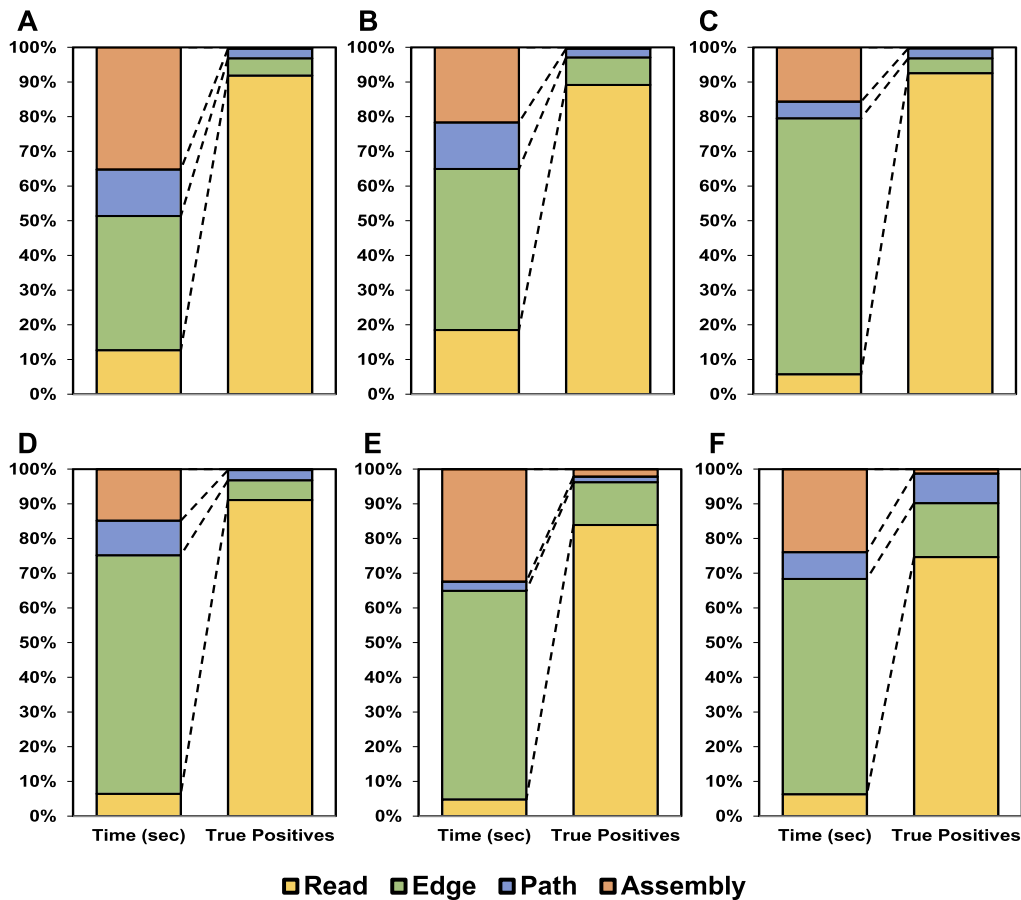


Figure 8. The breakdown of the running time – true prediction relation for the four main modules involved in iMPP(FGS). Yellow: application of FGS on reads; Green: assembly graph generation and the application of FGS on assembly graph edges; Blue: candidate path generation and the application of FGS on candidate paths; Orange: peptide assembly refinement. (A) DS1, (B) DS2, (C) DS3, (D) DS4, (E) DS5 and (F) DS6.

the local osmotic shock induced by the disease, and are suggested as predictive biomarkers for Type 2 diabetes (81).

DISCUSSION

In this work, we present a *de novo* metagenomic functional analysis workflow iMPP. iMPP directly operates on unassembled raw reads and is capable of discovering novel proteins or protein families. To the best of our knowledge, iMPP is currently the only method that integrates nucleotide assembly, gene calling, and peptide assembly based on their informational connection and dependency (Figure 1). The integration appears to be successful based on benchmark results. For gene calling, iMPP significantly improves upon the state-of-the-art methods FragGeneScan and MetaProdigal with a 4–20% higher sensitivity (Figure 4). Notably, iMPP has achieved a near-perfect recall rate of >90% on all metagenomic datasets (DS1–5), and a recall rate of ~78% on the metatranscriptome dataset (DS6) at a specificity level of ~85%. The highly accurate iMPP gene calling results further benefit downstream *de novo* peptide assembly, generating more peptide contigs with higher specificity. The improvements made by iMPP on the up-

stream gene calling and assembly further led to the reconstruction of meaningful metabolic pathways that are critical for functional investigation.

Note that iMPP is a generic workflow that can integrate different third-party software, which is not limited to the ones that we have tested in this article. For example, we have shown that iMPP(FGS) and iMPP(MPD) can improve the performances for FGS and MPD alone, indicating that the improvement is likely due to information integration, rather than dependency on a specific software. In this case, it is desirable to have iMPP to support more third-party software. Recall that iMPP has three main modules: the gene calling module, the nucleotide assembly module, and the peptide assembly module (Figure 1). At the current development stage, we have included FGS and MPD as two alternatives to the gene calling module. In the future, we further plan to include IDBA-UD (82) and MetaHit (83), in addition to the existing software SGA and SPAdes, as alternatives to the nucleotide assembly module. We also plan include MetaPA (84) and SFA-SPA (48), together with PLASS, as alternatives to the peptide assembly module. This effort will enable the search of a software combination for optimal performance, and higher flexibility to satisfy project-specific needs.

The peptide assembly results point to two seemingly counterintuitive observations regarding *de novo* peptide assembly. First, a more specific input set does not necessarily lead to a more specific assembly. As shown in Figure 4, FragGeneScan had a slightly higher specificity than iMPP (1–4%, see details in Supplementary Table S1). However, iMPP assembly showed higher contig- and read-level specificity than FGS + PLASS assembly (Table 3). The reason could be that the more specific input set generated by FragGeneScan was less comprehensive, and missing the true protein-coding reads made the assembly graph more fragmentary. It further resulted in many ultra short peptide contigs, which were subsequently filtered out, reducing the true positive rate and specificity. Second, a more comprehensive input set does not necessarily lead to a more complete assembly. Because iMPP only accepted the predicted coding reads as its input, its input was less complete than PLASS, which accepted all reads. Surprisingly, the iMPP assembly was more comprehensive than the PLASS assembly (Table 2 and Figure 6). The reason could be that contaminants (false pseudo peptides or mispredicted ORF from noncoding reads) may overlap with other peptides by chance, generating more false connections in the assembly graph. The false connections may confound graph traversal and reduce true positive output. As such, a more refined input that contains exactly all the coding reads will likely result in the best assembly. While these observations were made from peptide assembly, we believe that they also apply to nucleotide assembly, as most assembly algorithms are similar. The observation may provide insights to improve *de novo* nucleotide and peptide assembly from a different perspective: refining the input.

With the novel protein families discovered from the above analysis, as well as the known protein families, we expect to develop an algorithm to improve *de novo* nucleotide assembly. This work will complete the last piece of missing information flow from Figure 1 (the gray broken arrow). While it is possible to improve the assembly of individual genes using protein family profiles as a guide (26,54), it is unknown by how much it can improve assembly at the genome level. The improvement observed on individual genes suggests that a guided assembly can help resolve branches in the assembly graph. We shall take advantage of it towards more accurate genome assembly. With this module, we will further develop an iterative version of iMPP following the information flow shown in Figure 1. While the iterative version could be uneconomical given the already high recall rate of the current iMPP version (80–97%, Figure 4), it is of theoretical interest to investigate the limit of gene calling directly from fragmented sequences.

To promote practical applications of iMPP, we expect to include an additional module for hypothetical protein annotation. Note that a significant proportion (60–70%) of peptide sequences assembled by iMPP from the complete datasets cannot be aligned to the UniProt database (Table 5). Given the high contig-level specificity (~80%) observed from the subsampled datasets (Table 3), most of the assembled peptide contigs likely correspond to true novel proteins. Note that all of these assembled peptides are ≥ 60 aa, therefore they should contain sufficient information for reliable functional prediction. Specifically, we will develop a

hypothetical protein annotation module (85) that includes physicochemical property characterization, domain analysis, protein subcellular localization analysis, and protein–protein interaction analysis. We will also include a *de novo* clustering module to identify novel protein families and sequence motifs (86). Finally, we will also provide the corresponding DNA sequences of these proteins to facilitate their taxonomic analyses and experimental validations.

We also plan to improve the usability and efficiency of iMPP from a software engineering perspective. Currently, the iMPP workflow is wrapped with Nextflow (87), which ensures its reproducibility and facilitates its execution under different software environments. We will further leverage the power of Nextflow to modulate different components of iMPP to meet flexible needs in performance and efficiency. For example, the user will be able to eliminate the peptide assembly refinement step for a speedup without losing a significant number of true positive predictions. We will also try to speed up iMPP by ‘internalizing’ third-party software modules as libraries of iMPP. For example, iMPP first writes the assembly overlap graph and de Bruijn graph contigs into the hard disk and loads them to generate the hybrid graph. Internalizing the assembly graph generation modules can eliminate the hard-disk traffic and make the entire workflow more efficient. In addition, internalizing peptide assembly will also allow us to access the peptide assembly graph generated by the first PLASS run (for gene calling refinement); the information may help to save a significant amount of time for the second pass (for peptide contig reconstruction).

In conclusion, we present a novel method called iMPP for metagenomic functional analysis. iMPP integrates *de novo* nucleotide assembly, gene calling, and peptide assembly. iMPP is able to improve both gene calling and peptide assembly and has the potential to improve our current understanding of the functions of microbial communities. iMPP was implemented using GNU C++ and Python. It is freely available from (<https://github.com/Sirisha-t/iMPP>) under the Creative Commons BY-NC licence.

DATA AVAILABILITY

The data underlying this article are available in Figshare at <https://doi.org/10.6084/m9.figshare.22114178.v1>. The source code underlying this article is available in Figshare at <https://doi.org/10.6084/m9.figshare.22111187>.

SUPPLEMENTARY DATA

Supplementary Data are available at NARGAB Online.

FUNDING

ST, BL, AP and CZ are funded by National Science Foundation CAREER award [DBI-1943291]. Funding for open access charge: NSF CAREER award DBI-1943291.

Conflict of interest statement. None declared.

REFERENCES

- Chaparro, J.M., Sheflin, A.M., Manter, D.K. and Vivanco, J.M. (2012) Manipulating the soil microbiome to increase soil health and plant fertility. *Biol. Fertil. Soils*, **48**, 489–499.

2. Hamdan, L.J., Coffin, R.B., Sikaroodi, M., Greinert, J., Treude, T. and Gillevet, P.M. (2013) Ocean currents shape the microbiome of Arctic marine sediments. *ISME J*, **7**, 685–696.
3. Gusareva, E.S., Acerbi, E., Lau, K.J.X., Luhung, I., Premkrishnan, B.N.V., Kolundzija, S., Purbojati, R.W., Wong, A., Houghton, J.N.I., Miller, D. *et al.* (2019) Microbial communities in the tropical air ecosystem follow a precise diel cycle. *Proc. Natl. Acad. Sci. U.S.A.*, **116**, 23299–23308.
4. Khanna, S. and Tosh, P.K. (2014) A clinician's primer on the role of the microbiome in human health and disease. *Mayo Clin. Proc.*, **89**, 107–114.
5. Hollister, E.B., Gao, C. and Versalovic, J. (2014) Compositional and functional features of the gastrointestinal microbiome and their effects on human health. *Gastroenterology*, **146**, 1449–1458.
6. Valdes, A.M., Walter, J., Segal, E. and Spector, T.D. (2018) Role of the gut microbiota in nutrition and health. *BMJ*, **361**, k2179.
7. Cani, P.D., Bibiloni, R., Knauf, C., Waeg, A., Neyrinck, A.M., Delzenne, N.M. and Burcelin, R. (2008) Changes in gut microbiota control metabolic endotoxemia-induced inflammation in high-fat diet-induced obesity and diabetes in mice. *Diabetes*, **57**, 1470–1481.
8. Liu, R., Hong, J., Xu, X., Feng, Q., Zhang, D., Gu, Y., Shi, J., Zhao, S., Liu, W. and Wang, X. (2017) Gut microbiome and serum metabolome alterations in obesity and after weight-loss intervention. *Nat. Med.*, **23**, 859.
9. John, G.K. and Mullin, G.E. (2016) The gut microbiome and obesity. *Curr. Oncol. Rep.*, **18**, 45.
10. Koleva, P.T., Bridgman, S.L. and Kozyrskyj, A.L. (2015) The infant gut microbiome: evidence for obesity risk and dietary intervention. *Nutrients*, **7**, 2237–2260.
11. Barlow, G.M., Yu, A. and Mathur, R. (2015) Role of the gut microbiome in obesity and diabetes mellitus. *Nutr. Clin. Pract.*, **30**, 787–797.
12. Vatanen, T., Franzosa, E.A., Schwager, R., Tripathi, S., Arthur, T.D., Vehik, K., Lernmark, Å., Hagopian, W.A., Rewers, M.J. and She, J.-X. (2018) The human gut microbiome in early-onset type 1 diabetes from the TEDDY study. *Nature*, **562**, 589–594.
13. Halfvarson, J., Brislawn, C.J., Lamendella, R., Vázquez-Baeza, Y., Walters, W.A., Bramer, L.M., D'Amato, M., Bonfiglio, F., McDonald, D. and Gonzalez, A. (2017) Dynamics of the human gut microbiome in inflammatory bowel disease. *Nat. Microbiol.*, **2**, 17004.
14. Ananthkrishnan, A.N., Luo, C., Yajnik, V., Khalili, H., Garber, J.J., Stevens, B.W., Cleland, T. and Xavier, R.J. (2017) Gut microbiome function predicts response to anti-integrin biologic therapy in inflammatory bowel diseases. *Cell Host Microbe*, **21**, 603–610.
15. Franzosa, E.A., Sirota-Madi, A., Avila-Pacheco, J., Fornelos, N., Haiser, H.J., Reinker, S., Vatanen, T., Hall, A.B., Mallick, H. and McIver, L.J. (2019) Gut microbiome structure and metabolic activity in inflammatory bowel disease. *Nat. Microbiol.*, **4**, 293–305.
16. Harimoto, T. and Danino, T. (2019) Engineering bacteria for cancer therapy. *Emerg. Top. Life Sci.*, **3**, 623–629.
17. Nelson, M.H., Diven, M.A., Huff, L.W. and Paulos, C.M. (2015) Harnessing the microbiome to enhance cancer immunotherapy. *J. Immunol. Res.*, **2015**, 368736.
18. Nguyen, V.H. and Min, J.-J. (2017) Salmonella-mediated cancer therapy: roles and potential. *Nucl. Med. Mol. Imaging*, **51**, 118–126.
19. Huang, C., Li, M., Liu, B., Zhu, H., Dai, Q., Fan, X., Mehta, K., Huang, C., Neupane, P., Wang, F. *et al.* (2021) Relating gut microbiome and its modulating factors to immunotherapy in solid tumors: a systematic review. *Front. Oncol.*, **11**, 642110.
20. Chau, J., Yadav, M., Liu, B., Furqan, M., Dai, Q., Shahi, S., Gupta, A., Mercer, K.N., Eastman, E., Hejleh, T.A. *et al.* (2021) Prospective correlation between the patient microbiome with response to and development of immune-mediated adverse effects to immunotherapy in lung cancer. *BMC Cancer*, **21**, 808.
21. Venter, J.C., Remington, K., Heidelberg, J.F., Halpern, A.L., Rusch, D., Eisen, J.A., Wu, D., Paulsen, I., Nelson, K.E., Nelson, W. *et al.* (2004) Environmental genome shotgun sequencing of the Sargasso Sea. *Science*, **304**, 66–74.
22. Williamson, S.J. and Yooseph, S. (2012) From bacterial to microbial ecosystems (metagenomics). *Methods Mol. Biol.*, **804**, 35–55.
23. Segata, N., Waldron, L., Ballarini, A., Narasimhan, V., Jousson, O. and Huttenhower, C. (2012) Metagenomic microbial community profiling using unique clade-specific marker genes. *Nat. Methods*, **9**, 811–814.
24. Wood, D.E. and Salzberg, S.L. (2014) Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol.*, **15**, R46.
25. Nguyen, N.P., Mirarab, S., Liu, B., Pop, M. and Warnow, T. (2014) TIPP: taxonomic identification and phylogenetic profiling. *Bioinformatics*, **30**, 3548–3555.
26. Zhong, C., Edlund, A., Yang, Y., McLean, J.S. and Yooseph, S. (2016) Metagenome and metatranscriptome analyses using protein Family profiles. *PLoS Comput. Biol.*, **12**, e1004991.
27. Meinicke, P. (2015) UProC: tools for ultra-fast protein domain classification. *Bioinformatics*, **31**, 1382–1388.
28. Eddy, S.R. (2011) Accelerated profile HMM searches. *PLoS Comput. Biol.*, **7**, e1002195.
29. Suarez Araujo, C.P., Garcia Baez, P., Sanchez Rodriguez, A. and Santana Rodriguez, J.J. (2009) HUMANN-based system to identify benzimidazole fungicides using multi-synchronous fluorescence spectra: an ensemble approach. *Anal. Bioanal. Chem.*, **394**, 1059–1072.
30. Abubucker, S., Segata, N., Goll, J., Schubert, A.M., Izard, J., Cantarel, B.L., Rodriguez-Mueller, B., Zucker, J., Thiagarajan, M., Henrissat, B. *et al.* (2012) Metabolic reconstruction for metagenomic data and its application to the human microbiome. *PLoS Comput. Biol.*, **8**, e1002358.
31. The UniProt Consortium (2018) UniProt: the universal protein knowledgebase. *Nucleic Acids Res.*, **46**, 2699.
32. Altschul, S.F., Madden, T.L., Schaffer, A.A., Zhang, J., Zhang, Z., Miller, W. and Lipman, D.J. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.*, **25**, 3389–3402.
33. Peng, Y., Leung, H.C., Yiu, S.M. and Chin, F.Y. (2011) Meta-IDBA: a de novo assembler for metagenomic data. *Bioinformatics*, **27**, i94–i101.
34. Feldman, S.S. and Horan, T.A. (2011) Collaboration in electronic medical evidence development: a case study of the Social Security Administration's MEGAHIT System. *Int. J. Med Inform.*, **80**, e127–e140.
35. Namiki, T., Hachiya, T., Tanaka, H. and Sakakibara, Y. (2012) MetaVelvet: an extension of Velvet assembler to de novo metagenome assembly from short sequence reads. *Nucleic Acids Res.*, **40**, e155.
36. Afiahayati, Sato, K. and Sakakibara, Y. (2015) MetaVelvet-SL: an extension of the Velvet assembler to a de novo metagenomic assembler utilizing supervised learning. *DNA Res.*, **22**, 69–77.
37. Nurk, S., Meleshko, D., Korobeynikov, A. and Pevzner, P.A. (2017) metaSPAdes: a new versatile metagenomic assembler. *Genome Res.*, **27**, 824–834.
38. Delcher, A.L., Harmon, D., Kasif, S., White, O. and Salzberg, S.L. (1999) Improved microbial gene identification with GLIMMER. *Nucleic Acids Res.*, **27**, 4636–4641.
39. Lukashin, A.V. and Borodovsky, M. (1998) GeneMark. Hm: new solutions for gene finding. *Nucleic Acids Res.*, **26**, 1107–1115.
40. Hyatt, D., Chen, G.-L., Locascio, P.F., Land, M.L., Larimer, F.W. and Hauser, L.J. (2010) Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinf.*, **11**, 119–119.
41. Noguchi, H., Taniguchi, T. and Itoh, T. (2008) MetaGeneAnnotator: detecting species-specific patterns of ribosomal binding site for precise gene prediction in anonymous prokaryotic and phage genomes. *DNA Res.*, **15**, 387–396.
42. Rho, M., Tang, H. and Ye, Y. (2010) FragGeneScan: predicting genes in short and error-prone reads. *Nucleic Acids Res.*, **38**, e191.
43. Hoff, K.J., Lingner, T., Meinicke, P. and Tech, M. (2009) Orphelia: predicting genes in metagenomic sequencing reads. *Nucleic Acids Res.*, **37**, W101–W105.
44. Kelley, D.R., Liu, B., Delcher, A.L., Pop, M. and Salzberg, S.L. (2011) Gene prediction with Glimmer for metagenomic sequences augmented by classification and clustering. *Nucleic Acids Res.*, **40**, e9.
45. Zhu, W., Lomsadze, A. and Borodovsky, M. (2010) Ab initio gene identification in metagenomic sequences. *Nucleic Acids Res.*, **38**, e132.
46. Hyatt, D., LoCasio, P.F., Hauser, L.J. and Uberbacher, E.C. (2012) Gene and translation initiation site prediction in metagenomic sequences. *Bioinformatics*, **28**, 2223–2230.
47. Yang, Y. and Yooseph, S. (2013) SPA: a short peptide assembler for metagenomic data. *Nucleic Acids Res.*, **41**, e91.
48. Yang, Y., Zhong, C. and Yooseph, S. (2015) SFA-SPA: a suffix array based short peptide assembler for metagenomic data. *Bioinformatics*, **31**, 1833–1835.

49. Steinegger, M., Mirdita, M. and Soding, J. (2019) Protein-level assembly increases protein sequence recovery from metagenomic samples manyfold. *Nat. Methods*, **16**, 603–606.
50. Liu, J., Lian, Q., Chen, Y. and Qi, J. (2019) Amino acid based de Bruijn graph algorithm for identifying complete coding genes from metagenomic and metatranscriptomic short reads. *Nucleic Acids Res.*, **47**, e30.
51. Tang, H., Li, S. and Ye, Y. (2016) A graph-centric approach for metagenome-guided peptide and protein identification in metaproteomics. *PLoS Comput. Biol.*, **12**, e1005224.
52. Zhong, C., Yang, Y. and Yooseph, S. (2015) GRASP: guided reference-based assembly of short peptides. *Nucleic Acids Res.*, **43**, e18.
53. Huson, D.H., Tappu, R., Bazinet, A.L., Xie, C., Cummings, M.P., Nieselt, K. and Williams, R. (2017) Fast and simple protein-alignment-guided assembly of orthologous gene families from microbiome sequencing reads. *Microbiome*, **5**, 11.
54. Wang, Q., Fish, J.A., Gilman, M., Sun, Y., Brown, C.T., Tiedje, J.M. and Cole, J.R. (2015) Xander: employing a novel method for efficient gene-targeted metagenomic assembly. *Microbiome*, **3**, 32.
55. Simpson, J.T. and Durbin, R. (2012) Efficient de novo assembly of large genomes using compressed data structures. *Genome Res.*, **22**, 549–556.
56. Miller, J.R., Zhou, P., Mudge, J., Gurtowski, J., Lee, H., Ramaraj, T., Walenz, B.P., Liu, J., Stupar, R.M., Denny, R. *et al.* (2017) Hybrid assembly with long and short reads improves discovery of gene family expansions. *Bmc Genomics (Electronic Resource)*, **18**, 541.
57. Haghshenas, E., Asghari, H., Stoye, J., Chauve, C. and Hach, F. (2020) HASLR: fast hybrid assembly of long reads. *IScience*, **23**, 101389.
58. Antipov, D., Korobeynikov, A., McLean, J.S. and Pevzner, P.A. (2016) hybridSPAdes: an algorithm for hybrid assembly of short and long reads. *Bioinformatics*, **32**, 1009–1015.
59. Liu, B., Thippabhotla, S., Zhang, J. and Zhong, C. (2021) DRAGoM: classification and quantification of noncoding RNA in metagenomic data. *Front. Genet.*, **12**, 669495.
60. Zhong, C., Yang, Y. and Yooseph, S. (2017) *2017 IEEE 7th International Conference on Computational Advances in Bio and Medical Sciences (ICCABS)*, pp. 1.
61. Qin, J., Li, Y., Cai, Z., Li, S., Zhu, J., Zhang, F., Liang, S., Zhang, W., Guan, Y., Shen, D. *et al.* (2012) A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature*, **490**, 55–60.
62. Howe, A.C., Jansson, J.K., Malfatti, S.A., Tringe, S.G., Tiedje, J.M. and Brown, C.T. (2014) Tackling soil diversity with the assembly of large, complex metagenomes. *Proc. Natl. Acad. Sci. U.S.A.*, **111**, 4904–4909.
63. Biller, S.J., Berube, P.M., Dooley, K., Williams, M., Satinsky, B.M., Hackl, T., Hogle, S.L., Coe, A., Bergauer, K., Bouman, H.A. *et al.* (2018) Marine microbial metagenomes sampled across space and time. *Sci Data*, **5**, 180176.
64. Stewart, R.D., Auffret, M.D., Warr, A., Wisner, A.H., Press, M.O., Langford, K.W., Liachko, I., Snelling, T.J., Dewhurst, R.J., Walker, A.W. *et al.* (2018) Assembly of 913 microbial genomes from metagenomic sequencing of the cow rumen. *Nat. Commun.*, **9**, 870.
65. Sczyrba, A., Hofmann, P., Belmann, P., Koslicki, D., Janssen, S., Droge, J., Gregor, I., Majda, S., Fiedler, J., Dahms, E. *et al.* (2017) Critical Assessment of Metagenome Interpretation—a benchmark of metagenomics software. *Nat. Methods*, **14**, 1063–1071.
66. O’Leary, N.A., Wright, M.W., Brister, J.R., Ciuffo, S., Haddad, D., McVeigh, R., Rajput, B., Robbertse, B., Smith-White, B., Ako-Adjei, D. *et al.* (2016) Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.*, **44**, D733–D745.
67. Buchfink, B., Xie, C. and Huson, D.H. (2015) Fast and sensitive protein alignment using DIAMOND. *Nat. Methods*, **12**, 59–60.
68. Caspi, R., Billington, R., Ferrer, L., Foerster, H., Fulcher, C.A., Keseler, I.M., Kothari, A., Krummenacker, M., Latendresse, M., Mueller, L.A. *et al.* (2016) The MetaCyc database of metabolic pathways and enzymes and the BioCyc collection of pathway/genome databases. *Nucleic Acids Res.*, **44**, D471–D480.
69. Ye, Y. and Doak, T.G. (2009) A parsimony approach to biological pathway reconstruction/inference for genomes and metagenomes. *PLoS Comput. Biol.*, **5**, e1000465.
70. Kushkevych, I., Cejnar, J., Treml, J., Dordevic, D., Kollar, P. and Vitezova, M. (2020) Recent advances in metabolic pathways of sulfate reduction in intestinal bacteria. *Cells*, **9**, 698.
71. Kertesz, M.A. (2000) Riding the sulfur cycle—metabolism of sulfonates and sulfate esters in gram-negative bacteria. *FEMS Microbiol. Rev.*, **24**, 135–175.
72. Garcia, B., Merayo-Llodes, J., Rodriguez, D., Alcalde, I., Garcia-Suarez, O., Alfonso, J.F., Baamonde, B., Fernandez-Vega, A., Vazquez, F. and Quiros, L.M. (2016) Different use of cell surface glycosaminoglycans As adherence receptors to corneal cells by gram positive and gram negative pathogens. *Front. Cell Infect. Microbiol.*, **6**, 173.
73. Igarashi, K. and Kashiwagi, K. (2010) Modulation of cellular function by polyamines. *Int. J. Biochem. Cell Biol.*, **42**, 39–51.
74. Tofalo, R., Cocchi, S. and Suzzi, G. (2019) Polyamines and gut microbiota. *Front. Nutr.*, **6**, 16.
75. Matsumoto, M. and Benno, Y. (2007) The relationship between microbiota and polyamine concentration in the human intestine: a pilot study. *Microbiol. Immunol.*, **51**, 25–35.
76. Noack, J., Dongowski, G., Hartmann, L. and Blaut, M. (2000) The human gut bacteria *Bacteroides thetaiotaomicron* and *Fusobacterium varium* produce putrescine and spermidine in cecum of pectin-fed gnotobiotic rats. *J. Nutr.*, **130**, 1225–1231.
77. Craig, S.A. (2004) Betaine in human nutrition. *Am. J. Clin. Nutr.*, **80**, 539–549.
78. Koistinen, V.M., Karkkainen, O., Borewicz, K., Zarei, I., Jokkala, J., Micard, V., Rosa-Sibakov, N., Auriola, S., Aura, A.M., Smidt, H. *et al.* (2019) Contribution of gut microbiota to metabolism of dietary glycine betaine in mice and in vitro colonic fermentation. *Microbiome*, **7**, 103.
79. Richter, A.A., Mais, C.N., Czech, L., Geyer, K., Hoepfner, A., Smits, S.H.J., Erb, T.J., Bange, G. and Bremer, E. (2019) Biosynthesis of the stress-protectant and chemical chaperon ectoine: biochemistry of the transaminase EctB. *Front. Microbiol.*, **10**, 2811.
80. Larsen, P.E. and Dai, Y. (2015) Metabolome of human gut microbiome is predictive of host dysbiosis. *Gigascience*, **4**, 42.
81. Shen, N., Dimitrova, N., Ho, C.H., Torres, P.J., Camacho, F.R., Cai, Y., Vuyisich, M., Tanton, D. and Banavar, G. (2021) Gut microbiome activity predicts risk of type 2 diabetes and metformin control in a large human cohort. medRxiv doi: <https://doi.org/10.1101/2021.08.13.21262051>, 18 August 2021, preprint: not peer reviewed.
82. Peng, Y., Leung, H.C., Yiu, S.M. and Chin, F.Y. (2012) IDBA-UD: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics*, **28**, 1420–1428.
83. Li, D., Liu, C.M., Luo, R., Sadakane, K. and Lam, T.W. (2015) MEGAHIT: an ultra-fast single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics*, **31**, 1674–1676.
84. Liu, J., Lian, Q., Chen, Y. and Qi, J. (2019) Amino acid based de Bruijn graph algorithm for identifying complete coding genes from metagenomic and metatranscriptomic short reads. *Nucleic Acids Res.*, **47**, e30.
85. Ijaq, J., Chandrasekharan, M., Poddar, R., Bethi, N. and Sundararajan, V.S. (2015) Annotation and curation of uncharacterized proteins - challenges. *Front Genet*, **6**, 119.
86. Bailey, T.L., Boden, M., Buske, F.A., Frith, M., Grant, C.E., Clementi, L., Ren, J., Li, W.W. and Noble, W.S. (2009) MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res.*, **37**, W202–W208.
87. Di Tommaso, P., Chatzou, M., Floden, E.W., Barja, P.P., Palumbo, E. and Notredame, C. (2017) Nextflow enables reproducible computational workflows. *Nat. Biotechnol.*, **35**, 316–319.