

Scope is all you need:

Transforming LLMs for HPC Code



Contributors

Tal Kadosh

Ben-Gurion University, IAEC, Israel

Niranjan Hasabnis

Intel Labs, United States

Vy A. Vo

Intel Labs, United States

Nadav Schneider

Ben-Gurion University, IAEC, Israel

Neva Krien

Independent Researcher, Israel

Abdul Wasay

Intel Labs, United States

Nesreen Ahmed

Intel Labs, United States

Ted Willke

Intel Labs, United States

Guy Tamir

Intel, United States

Yuval Pinter

Ben-Gurion University, Israel

Timothy Mattson

Intel Labs, United States

Gal Oren

Technion, NRCN, Israel

galoren@cs.technion.ac.il

and Sponsors



Centers of Excellence



**Hewlett Packard
Enterprise**



CRAY

Our 2023 papers in this topic (so far)

- Kadosh, Tal, Niranjan Hasabnis, Vy A. Vo, Nadav Schneider, Neva Krien, Abdul Wasay, Nesreen Ahmed, Ted Willke, Guy Tamir, Yuval Pinter, Timothy Mattson, and Gal Oren . "**Scope is all you need: Transforming LLMs for HPC Code.**" *arXiv preprint arXiv:2308.09440* (2023).
- Kadosh, Tal, Niranjan Hasabnis, Timothy Mattson, Yuval Pinter, and Gal Oren. "**Quantifying OpenMP: Statistical Insights into Usage and Adoption [HPCorpus].**" *arXiv preprint arXiv:2308.08002* (2023). [HPEC'23](#)
- Schneider, Nadav, Tal Kadosh, Niranjan Hasabnis, Timothy Mattson, Yuval Pinter, and Gal Oren. "**MPI-rical: Data-Driven MPI Distributed Parallelism Assistance with Transformers.**" *arXiv preprint arXiv:2305.09438* (2023). [AI4DEV @ SC'23](#)
- Kadosh, Tal, Nadav Schneider, Niranjan Hasabnis, Timothy Mattson, Yuval Pinter, and Gal Oren. "**Advising OpenMP Parallelization via a Graph-Based Approach with Transformers.**" *arXiv preprint arXiv:2305.11999* (2023). [IWOMP'23](#)
- Harel, Re'em, Yuval Pinter, and Gal Oren. "**Learning to parallelize in a shared-memory environment with transformers.**" *Proceedings of the 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming*. 2023. [PPoPP'23](#)
 - Re'em Harel, Tal Kadosh, Niranjan Hasabnis, Timothy Mattson, Yuval Pinter, and Gal Oren. "**PragFormer: Data-driven Parallel Source Code Classification with Transformers**", 29 August 2023, *preprint* (V1) Research Square [<https://doi.org/10.21203/rs.3.rs-3254961/v1>]

Scope is all you need:

Intro and objectives - LLMs for HPC:

Gal Oren

Novel HPC code database - HPCorpus:

Niranjan Hasabnis

Compiler-oriented tokenizer - Tokompiler:

Vy A. Vo

Downstream Task #1 - OpenMP:

Tal Kadosh

Downstream Task #2 - MPI:

Nadav Schneider

Scope is all you need:

Intro and objectives - LLMs for HPC:

Gal Oren

Novel HPC code database - HPCorpus:

Niranjan Hasabnis

Compiler-oriented tokenizer - Tokompiler:

Vy A. Vo

Downstream Task #1 - OpenMP:

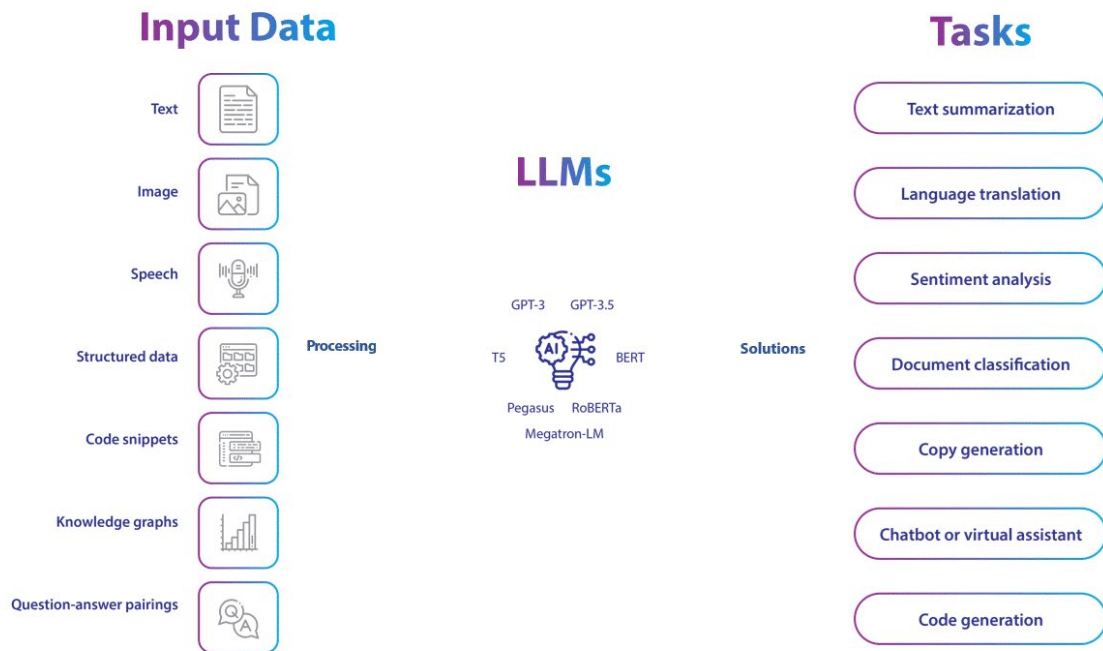
Tal Kadosh

Downstream Task #2 - MPI:

Nadav Schneider

Ever-increasing scope

- Infinite data
- Extend data representations
- (Really) large models
- Huge compute power
- Increasing tasks
- Fusion of knowledge
- ...



General Intelligence

Sparks of Artificial General Intelligence: Early experiments with GPT-4

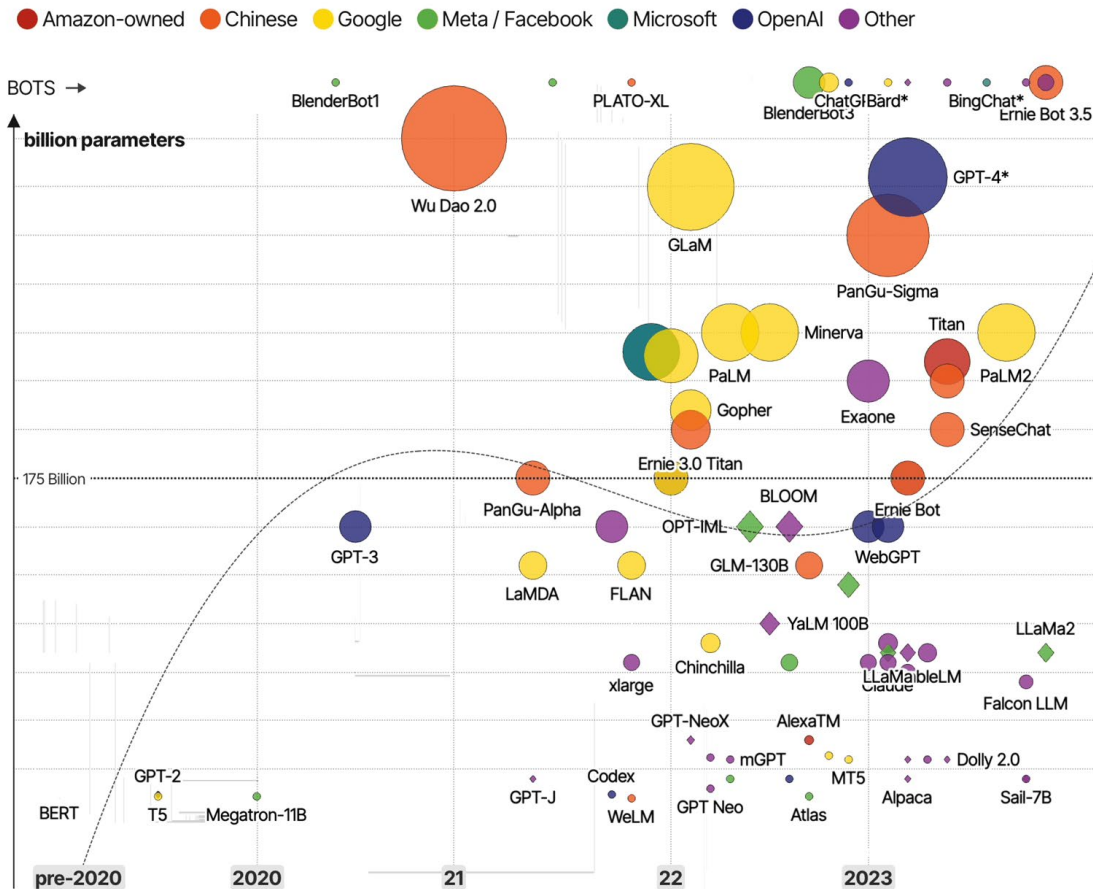
Sébastien Bubeck Varun Chandrasekaran Ronen Eldan Johannes Gehrmann
 Eric Horvitz Ece Kamar Peter Lee Yin Tat Lee Yuanzhi Li Scott Lundberg
 Harsha Nori Hamid Palangi Marco Tulio Ribeiro Yi Zhang

Microsoft Research

Abstract

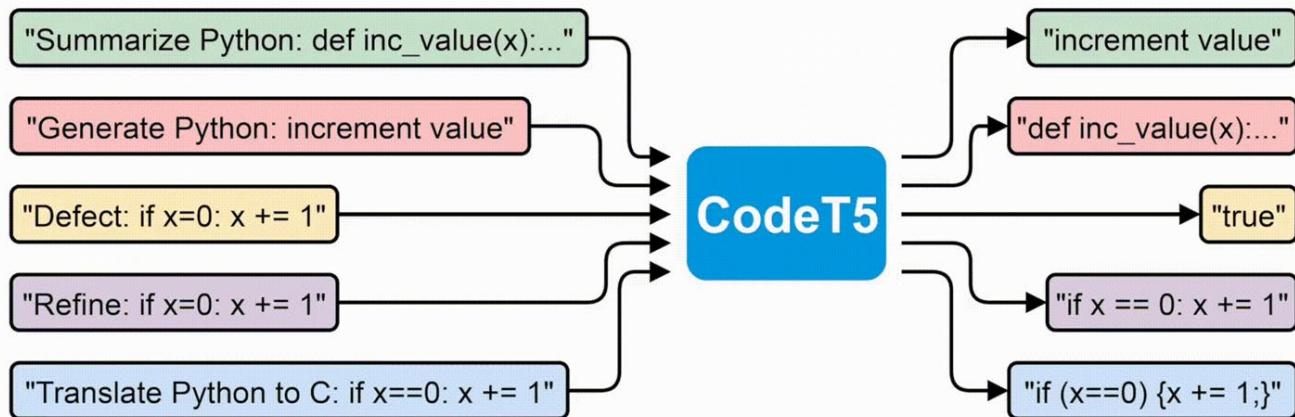
Artificial intelligence (AI) researchers have been developing and refining large language models (LLMs) that exhibit remarkable capabilities across a variety of domains and tasks, challenging our understanding of learning and cognition. The latest model developed by OpenAI, GPT-4 [Ope23], was trained using an unprecedented scale of compute and data. In this paper, we report on our investigation of an early version of GPT-4, when it was still in active development by OpenAI. We contend that (this early version of) GPT-4 is part of a new cohort of LLMs (along with ChatGPT and Google's PaLM for example) that exhibit more general intelligence than previous AI models. We discuss the rising capabilities and implications of these models. We demonstrate that, beyond its mastery of language, GPT-4 can solve novel and difficult tasks that span mathematics, coding, vision, medicine, law, psychology and more, without needing any special prompting. Moreover, in all of these tasks, GPT-4's performance is strikingly close to human-level performance, and often vastly surpasses prior models such as ChatGPT. Given the breadth and depth of GPT-4's capabilities, we believe that it could reasonably be viewed as an early (yet still incomplete) version of an artificial general intelligence (AGI) system. In our exploration of GPT-4, we put special emphasis on discovering its limitations, and we discuss the challenges ahead for advancing towards deeper and more comprehensive versions of AGI, including the possible need for pursuing a new paradigm that moves beyond next-word prediction. We conclude with reflections on societal influences of the recent technological leap and future research directions.

Source: <https://arxiv.org/abs/2303.12712>



Source: <https://informationisbeautiful.net/visualizations/the-rise-of-generative-ai-large-language-models-llm-s-like-chatgpt/>

Code Intelligence



Code Generation
Code Debugging
Code Optimization
Code Explanation
Algorithm Design
Data Manipulation
Web Development
DB Interaction
Code Review
API Integration
Code Documentation
Code Translation
Code Testing
Security Analysis

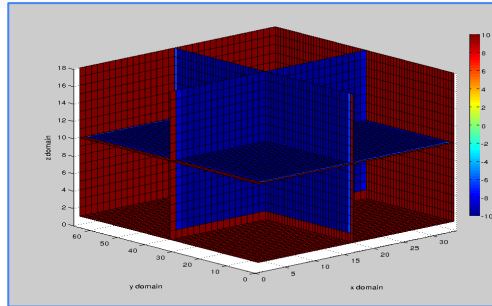
HPC Code - Same objectives?

HPC Langs

```
example.F90
1  ! Test program
2  PROGRAM example
3  USE test_module
4  IMPLICIT NONE
5  REAL :: a
6  TYPE(my_type) :: ex_var
7
8  a=4
9  CALL
10
11 END PROGRAM
12
```

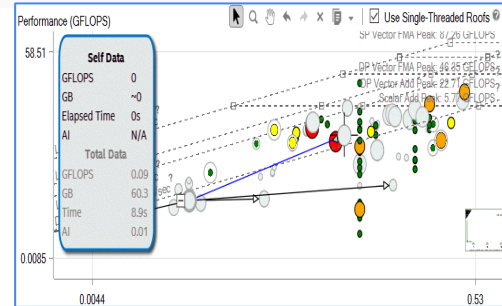
Source: <https://github.com/hansce/autocomp-lete-fortran>

Scientific Comp



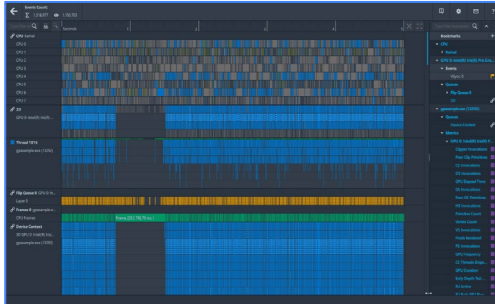
Source: https://dournac.org/info/parallel_heat3d

Performance



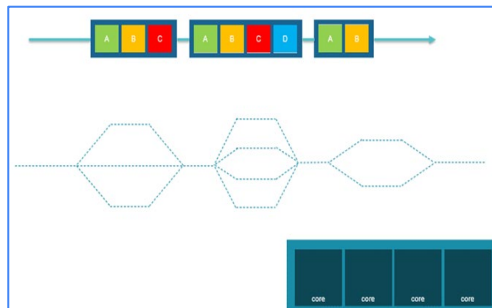
Source: intel.com/content/www/..../roffline-with-callstacks/

Dependencies



Source: intel.com/content/www/..../graphics-performance...html

Hardware Arch



Source: <http://ferestrepoa.github.io/paradigmas-de-programacion>

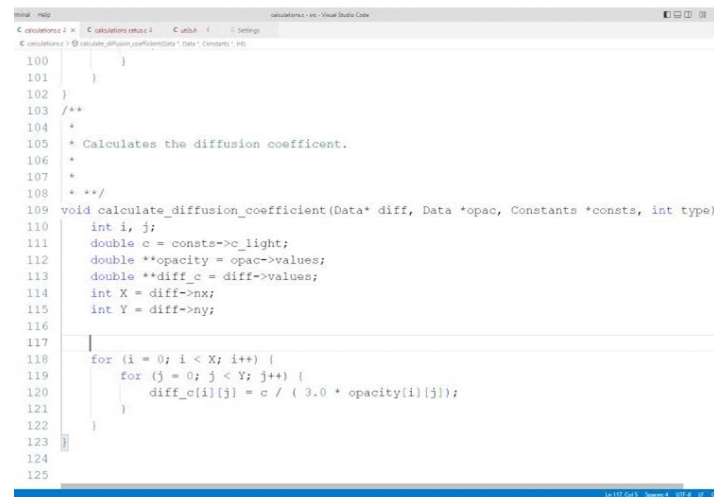
Parallelization

```
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

Source: <https://www.cs.technion.ac.il/view-event.php?evId=10334>

HPC Code Intelligence

- Speed up code development
- Generate **optimized parallel** code from high-level specifications
- Improve code readability and maintainability as **legacy code**
- Facilitate **cross-platform** compatibility through code translation
- Optimize code for **parallel execution** and **resource utilization**
- Identify **performance bottlenecks** and parallelization issues
- Enable code search for existing **optimized implementations**
- Facilitate **correct API** usage
- **Analyze code dependencies**
- Enforce **performance-oriented coding** standards

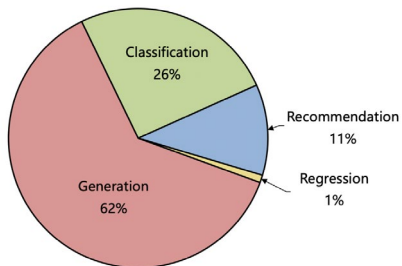
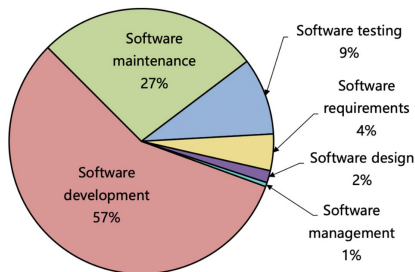


```
100     }
101   }
102 }
103 /**
104  *
105  * Calculates the diffusion coefficient.
106  *
107  *
108  */
109 void calculate_diffusion_coefficient(Data* diff, Data *opac, Constants *consts, int type)
110 {
111     int i, j;
112     double c = consts->c_light;
113     double **opacity = opac->values;
114     double **diff_c = diff->values;
115     int X = diff->nX;
116     int Y = diff->nY;
117
118     for (i = 0; i < X; i++) {
119         for (j = 0; j < Y; j++) {
120             diff_c[i][j] = c / (3.0 * opacity[i][j]);
121         }
122     }
123 }
124
125
```

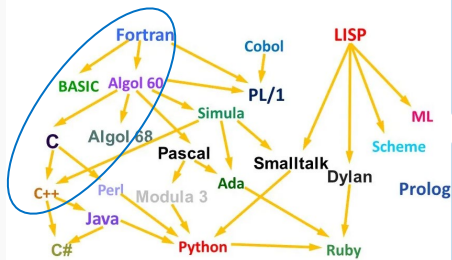
Suggest shared-memory parallelization on-the-fly based on LLM

What is our baseline? (1/3)

- Can non-HPC languages be the baseline for HPC tasks?
- Are there any relevant large corpora of high-performance code?
- What are the exact downstream tasks for HPC, how to correlate those with common models?
- Is mixing of languages and tasks is even good for HPC LLMs?



Hou, Xinyi, et al. "Large Language Models for Software Engineering: A Systematic Literature Review." arXiv preprint arXiv:2308.10620 (2023).



Source: <https://medium.com/@anaharris/human-languages-vs-programming-languages-c89410f13252>

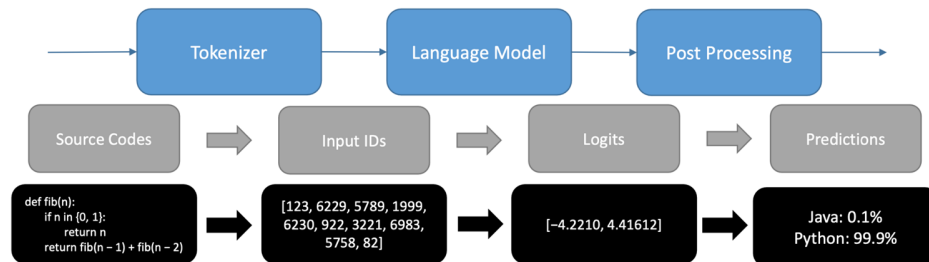
Framework	Year	Task(s)	Baseline(s)	Supported Languages(s)	Open Sourced
Refactory [137]	2019	Defect Detection	BLEU	Java	✗
CuBERT [138]	2020	Code Refinement, Defect Detection	BERT	Python	✓
CuGLM [139]	2020	Code Completion	BERT	Java, TypeScript	✓
Intellicode [140]	2020	Code Generation, Code Completion	GPT-2	Python, C#, JavaScript, and TypeScript	✗
Great [141]	2020	Defect Detection	Transformers	Python	✓
TreeGEN [51]	2020	Code Generation	Vanilla	Python	✓
C-BERT [127]	2020	Defect Detection	Transformers	BERT	✗
TransCoder [142]	2020	Code Translation	Vanilla	C++, Java, and Python	✗
GraphCodeBERT [143]	2020	Code Summarization, Code Refinement	Transformers	Java	✗
Codex [35]	2021	Code Generation, Code Completion, Code Summarization, Benchmark	BERT	JavaScript, Go, Perl, and 6 more	✗
Copilot [144]	2021	Code Generation, Code Completion	GPT-3	Java, PHP, Python, and 5 more	✗
CodeT5 [145]	2021	Code Summarization, Code Generation, Code Translation, Code Refinement, Defect Detection, Clone Detection	Codex	Python, Java	✓
Tfix [146]	2021	Code Refinement, Defect Detection	T5	JavaScript	✓
CodeRL [147]	2021	Code Summarization, Code Generation, Code Translation, Code Refinement, Defect Detection, Clone Detection	T5	Java	✓
TreeBERT [148]	2021	Code Summarization	Vanilla	Python, Java	✓
BUGLAB [149]	2021	Code Refinement, Defect Detection	Transformers	Python	✓
TBCC [150]	2021	Clone Detection	GREAT	C, Java	✓
APPS [36]	2021	Benchmark	Transformers	Python	✓
CodeXGLUE [34]	2021	Benchmark	N/A	Python	✓
CoTeXT [151]	2021	Code Summarization, Code Generation, Code Refinement, Defect detection	N/A	Python, Java, JavaScript, PHP, Ruby, Go	✓
SynCoBERT [152]	2021	Code Translation, Defect Detection, Clone Detection	T5	Ruby, JavaScript, Go, Python, Java, PHP	✗
TravTrans [153]	2021	Code Completion	BERT	Python	✗
CCAG [154]	2021	Code Completion	Vanilla	JavaScript, Python	✗
DeepDebug [155]	2021	Defect Detection	Transformers	Java	✓
Recorder [93]	2021	Defect Detection	Reformer	Java	✓
PLBART [156]	2021	Code Summarization, Code Generation, Code Translation, Code Refinement, Clone Detection, Detect Detection	TreeGen	Java	✓
PLBART [156]	2021	Code Summarization, Code Generation, Code Translation, Code Refinement, Clone Detection, Detect Detection	BART	Java, Python	✗
CODEGEN [157]	2022	Code Generation	GPT-NEO & GPT	Python	✓
GPT-2 for APR [158]	2022	Code Refinement	GPT-2	JavaScript	✓
CERT [39]	2022	Code Generation	CODEGEN	Python	✓
PyCoder [87]	2022	Code Generation	GPT-2	Python	✓
AlphaCode [38]	2022	Code Generation	GPT	Java	✗
InCoder [40]	2022	Code Generation, Code Completion, Code Summarization	GPT-3	Java, JavaScript, Python	✓
RewardRepair [159]	2022	Code Refinement, Defect Detection	T5	Java	✓
CodeParrot [37]	2022	Code Refinement, Defect Detection	GPT-2	Python	✓
AlphaRepair [160]	2022	Code Refinement, Defect Detection	CodeBERT	Java	✓
CodeReviewer [128]	2022	Code Summarization, Code Refinement, Defect Detection	CodeT5	Java	✓
TransRepair [161]	2022	Code Refinement, Defect Detection	BLEU	Java	✗
NatGen [162]	2022	Code Generation, Code Translation, Code Refinement	CodeT5	Java, Python, Go, JavaScript, Ruby, PHP	✓
DuaISC [163]	2022	Code Generation, Code Summarization	T5	Shellcode	✓
VulRepair [164]	2022	Code Refinement, Defect Detection	T5	C, C++	✓
CoditT5 [165]	2022	Code Summarization, Defect Detection	CodeT5	Java, Python, Ruby, PHP, Go	✓
C4 [166]	2022	Clone Detection	CodeBERT	C++, C#, Java, Python	✓
SPT-Code [167]	2022	Code Summarization, Code Completion, Code Refinement, Code Translation	CodeBERT & GraphCodeBERT	Python, Java, JavaScript, PHP, Go	✓
ExploitGen [168]	2023	Code Generation	CodeBERT	Python, Assembly	✓
SantaCoder [169]	2023	Code Summarization, Code Generation	GPT-2	Python, Java, and JavaScript	✓
xCodeEval [42]	2023	Benchmark	N/A	Python, Java, C++, PHP, and 8 more	✓
StarCoder [170]	2023	Code Generation, Code Completion, Code Summarization	BERT & SantaCoder	HTML, Python, Java, and 83 more	✓

Wong, Man-Fai, et al. "Natural Language Generation and Understanding of Big Code for AI-Assisted Programming: A Review." Entropy 25.6 (2023): 888.

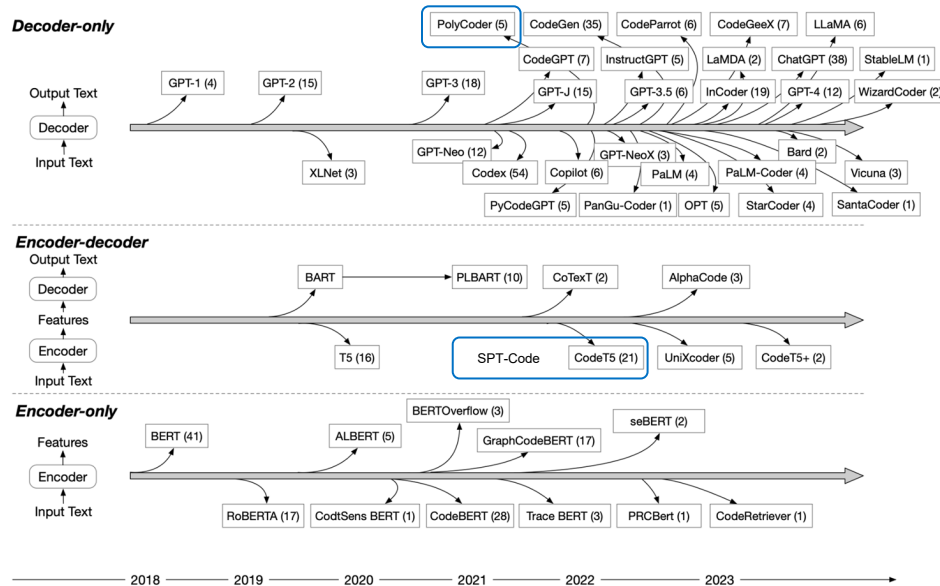
What is our baseline? (2/3)

- Are common tokenizers even relevant for compilation oriented code?
- What is the relevant model architecture for HPC tasks learning?
 - Encoder: Classification
 - Decoder: Generation
 - Encoder-Decoder: Both

(understanding input sequences and generating output sequences, often with different lengths and structures)



Wong, Man-Fai, et al. "Natural Language Generation and Understanding of Big Code for AI-Assisted Programming: A Review." Entropy 25.6 (2023): 888.



What is our baseline? (3/3)

- Should we consider other representations rather just plain code or AST?
- What about hardware specs? Telemetry? Should it be part of the input?
- Can we reduce the size of the input?

○ Of the model?

■ Of the training?

Code sample:

```
int main(){
    int r[2800 + 1];
    int i, k;
    int b, d;
    int c = 0;

    for(i = 0; i < 2800; i++){
        r[i] = 2000;
    }

    for(k = 2800; k > 0; k -= 14) {
        d = 0;
        i = k;

        for(;;){
            d += r[i] * 10000;
            b = 2 * i - 1;

            r[i] = d % b;
            d /= b;
            if(i == 0) break;
            d *= i;
        }
        printf("%4d", c + d / 10000);
        c = d % 10000;
    }

    return 0;
}
```

Tokenized code:

```
['int', 'main', '(', ')', '{', 'int', 'r', '[', '28', '00', '+', '1', ']', ',', 'int', 'i', ',', 'int', 'k', ',', 'int', 'b', ',', 'd', ',', 'int', 'c', '=', '0', ',', 'for', '(', 'i', '=', '0', ',', 'i', '<', '28', '00', ',', 'i', '++', ')', '{', 'r', '[', 'i', ']', '=', '2000', ',', '}', 'for', '(', 'k', '=', '28', '00', ',', 'k', '>', '0', ',', 'k', '-=', '14', ')', '{', 'd', '=', '0', ',', 'i', '=', 'k', ',', 'for', '(', ';;', ')', '{', 'd', '+=', 'r', '[', 'i', ']', '*', '10000', ',', 'b', '=', '2', '*', 'i', '-', '1', ',', 'r', '[', 'i', ']', '=', 'd', '%', 'b', ',', 'd', '/', 'b', ',', 'if', '(', 'i', '=', '0', ')', 'break', ',', 'd', '*', 'i', ',', '}', 'printf', '(', '"', '%.4d', '"', ',', 'c', '+', 'd', '/', '10000', ')', ',', 'c', '=', 'd', '%', '10000', ',', '}', 'return', '0', ',', '}' ]
```

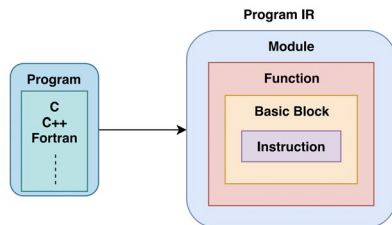
Tokenized AST:

```
['translation_unit', '(', 'function_definition', '(', 'int', 'function_declarator', '(', 'main', 'parameter_list', '(', '(', ')', ')', ')', 'compound_statement', '(', '(', 'declaration', '(', 'int', 'array_declarator', '(', 'r', '[', '2800', ']', ')', ')', 'compound_statement', '(', '2800', '+', '1', ')', ')', ')', 'declaration', '(', 'int', 'i', ',', 'k', ',', 'd', ',', 'int', 'b', ',', 'd', ',', 'd', ',', 'int', 'c', '=', '0', ')', 'declaration', '(', 'int', 'init_declarator', 'expression_statement', '(', 'assignment_expression', '(', 'subscript_expression', '(', 'r', '[', 'i', ']', ')', ')', ')', ... ]
```

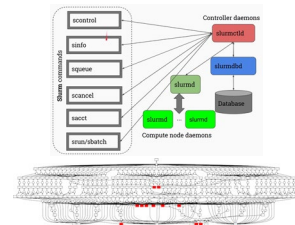
Tokenized DFG:

```
[('c', 24, 'comesFrom', [], []), ('i', 30, 'computedFrom', [0], [32]), ('0', 32, 'comesFrom', [], []), ('i', 34, 'comesFrom', [i], [30]), ('i', 38, 'comesFrom', [i], [30]), ('r', 42, 'computedFrom', [2000], [47]), ('i', 44, 'computedFrom', [2000], [47]), ('2000', 47, 'comesFrom', [i], []), ('k', 52, 'computedFrom', [2800], [54]), ('2800', 54, 'comesFrom', [], []), ('i', 56, 'comesFrom', [k], [52]), ('k', 60, 'computedFrom', [i], [62]), ('14', 62, 'comesFrom', [], []), ('d', 65, 'computedFrom', [0], [67]), ('0', 67, 'comesFrom', [], []), ... ]
```

IR ? loop-carried dependencies? hardware specs? Telemetry? ...



Terminal	Terminal
hostnamectl	free -h
uname -a	lspci
lscpu	lsusb
lshw (if installed)	ifconfig
lsblk	lsmod
df -h	ls /proc



“Scope is all you need”

(if you want transforming LLMs for HPC)

- The authors

Scope is all you need:

Intro and objectives - LLMs for HPC: Gal Oren

Novel HPC code database - HPCorpus: **Niranjan Hasabnis**

Compiler-oriented tokenizer - Tokompiler: Vy A. Vo

Downstream Task #1 - OpenMP: Tal Kadosh

Downstream Task #2 - MPI: Nadav Schneider

HPC languages

- Frontier and Aurora Supercomputers mostly rely on C, C++, Fortran...
- Top applications for HPC rely on MPI for distribution and OpenMP for shared-memory parallelization and GPU offloading
- So why not collect large corpus of those for LLM training, fine-tuning and special downstream tasks?

Domain	Application	Programming Languages, Models, Portability Layers	Algorithmic / Motifs
Astrophysics	Cholla	C/C++, HIP	Finite volume, multi-D FFTs
Molecular Dynamics	NAMD	C++, HIP, OpenMP, Charm++	Particle mesh Ewald, MD
Materials Science	LSMS	Fortran, C/C++, HIP, OpenMP	Dense linear algebra
Biology	CoMet	C++, HIP, OpenMP	Dense linear algebra
Fluid Dynamics	GESTS	Fortran, OpenMP	FFTs, Fourier pseudo-spectral
Nuclear Physics	NUCCOR	Fortran, C, OpenMP, HIP	Sparse and dense linear algebra, Hartree-Fock
Plasma Physics	PIConGPU	C++, HIP, OpenMP, CUPLA & ALPAKA	Particle in cells
Subsurface Flow	LBPM	C++, HIP	Lattice Boltzmann



Source: W. Joubert, "Ready for the Frontier" presentation, ISC'23



Aurora
Leadership Computing Facility
Exascale Supercomputer

Peak Performance
≥ 2 Exaflops DP

Intel GPU
Ponte Vecchio (PVC)

Intel Xeon Processor
Sapphire Rapids with High Bandwidth Memory

Platform
HPE Cray-Ex

Compute Node
2 Xeon SPR+HBM processors
6 Ponte Vecchio GPUs
Node Unified Memory Architecture
8 fabric endpoints

GPU Architecture
Intel XeHPC architecture
High Bandwidth Memory Stacks

Node Performance
>130 TF

System Size
>9,000 nodes

Aggregate System Memory
>10 PB aggregate System Memory

System Interconnect
HPE Slingshot 11
Dragonfly topology with adaptive routing

Network Switch
25.6 Tb/s per switch (64 200 Gb/s ports)
Links with 25 GB/s per direction

High-Performance Storage
220 PB
≥25 TB/s DAOS bandwidth

Software Environment

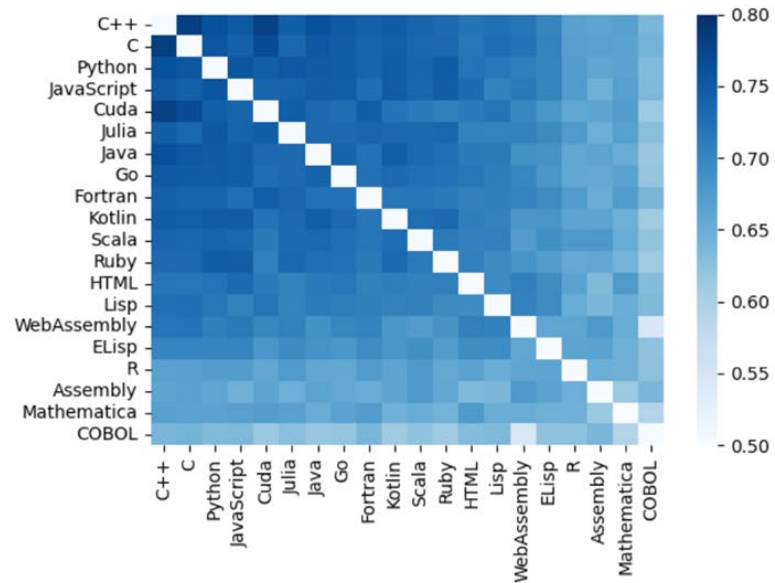
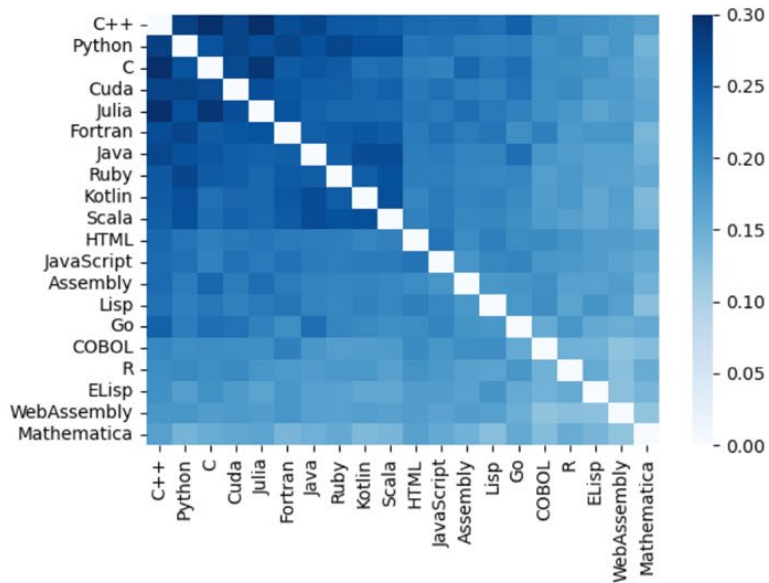
- C/C++
- Fortran
- SYCL/DPC++
- OpenMP offload
- Kokkos
- RAJA
- Intel Performance Tools

Software Environment

- C/C++
- Fortran
- SYCL/DPC++
- OpenMP offload
- Kokkos
- RAJA
- Intel Performance Tools

Source: ALCF SDL Workshop October 6, 2022, Colleen Bertoni and Scott Parker

Pairwise similarity between the representation of tokens of all code languages in CodeBERT



Non-pretrained setting

Pretrained setting

Source: Katzy, Jonathan, Maimen Izadi, and Arie van Deursen. "On the Impact of Language Selection for Training and Evaluating Programming Language Models." arXiv preprint arXiv:2308.13354 (2023).

HP Corpus

- All of C, C++ and Fortran codes from 2012-mid 2023
- High Level:
 - ~300K repos
 - ~70 GB
 - ~9 M files
 - ~155 M functions



```
1 WITH selected_repos as (  
2 SELECT f.id, f.repo_name as repo_name, f.ref as  
3   ↪ ref, f.path as path  
4 FROM `bigquery-public-data.github_repos.files` as  
5   ↪ f  
6 JOIN `bigquery-public-data.github_repos.licenses`  
7   ↪ as l on l.repo_name = f.repo_name  
8 ),  
9 deduped_files as (  
10 SELECT f.id, MIN(f.repo_name) as repo_name,  
11   ↪ MIN(f.ref) as ref, MIN(f.path) as path  
12 FROM selected_repos as f  
13 GROUP BY f.id  
14 )  
15 SELECT  
16   f.repo_name, f.ref, f.path, c.copies, c.content,  
17 FROM deduped_files as f  
18 JOIN `bigquery-public-data.github_repos.contents`  
19   ↪ as c on f.id = c.id  
20 WHERE  
21   NOT c.binary  
22   AND (f.path like '%.c' OR f.path like '%.cpp' OR  
23     ↪ f.path like '%.f' OR f.path like '%.f90' OR  
24     ↪ f.path like '%.f95')
```

	Repos	Size(GB)	Files (#)	Functions (#)
C	144,522	46.23	4,552,736	87,817,591
C++	150,481	26.16	4,735,196	68,233,984
Fortran	3,683	0.68	138,552	359,272

HP Corpus: statistical analysis

- Insights into parallel programming APIs and OpenMP
- Presented in our “Quantifying OpenMP” paper
- To be presented at IEEE High Performance Extreme Computing (HPEC) 2023

Quantifying OpenMP: Statistical Insights into Usage and Adoption

Tal Kadosh^{1,2}, Niranjan Hasabnis³, Timothy Mattson³, Yuval Pinter¹ and Gal Oren^{4,5}

¹Department of Computer Science, Ben-Gurion University, Israel

²Israel Atomic Energy Commission

³Intel Labs, United States

⁴Scientific Computing Center, Nuclear Research Center – Negev, Israel

⁵Department of Computer Science, Technion – Israel Institute of Technology, Israel

talkad@post.bgu.ac.il, niranjan.hasabnis@intel.com,

timothy.g.mattson@intel.com, pinter@bgu.ac.il, galoren@cs.technion.ac.il

Abstract—In high-performance computing (HPC), the demand for efficient parallel programming models has grown dramatically since the end of Dennard Scaling and the subsequent move to multi-core CPUs. OpenMP stands out as a popular choice due to its simplicity and portability, offering a directive-driven approach for shared-memory parallel programming. Despite its wide adoption, however, there is a lack of comprehensive data on the actual usage of OpenMP constructs, hindering unbiased insights into its popularity and evolution.

This paper presents a statistical analysis of OpenMP usage and adoption trends based on a novel and extensive database, HPCORPUS, compiled from GitHub repositories containing C, C++, and Fortran code. The results reveal that OpenMP is the dominant parallel programming model, accounting for 45% of all analyzed parallel APIs. Furthermore, it has demonstrated steady and continuous growth in popularity over the past decade. Analyzing specific OpenMP constructs, the study provides in-depth insights into their usage patterns and preferences across the three languages. Notably, we found that while OpenMP has a strong “common core” of constructs in common usage (while the rest of the API is less used), there are new adoption trends as well, such as `simd` and `target` directives for accelerated computing and `task` for irregular parallelism.

Overall, this study sheds light on OpenMP’s significance in HPC applications and provides valuable data for researchers and practitioners. It showcases OpenMP’s versatility, evolving adoption, and relevance in contemporary parallel programming, underlining its continued role in HPC applications and beyond. These statistical insights are essential for making informed decisions about parallelization strategies and provide a foundation for further advancements in parallel programming models and techniques.

HPCORPUS, as well as the analysis scripts and raw results, are available at: <https://github.com/Scientific-Computing-Lab-NRC/NHPCorpus>

Index Terms—HPCorpus, BigQuery, GitHub, C, C++, Fortran, OpenMP, MPI, OpenCL, CUDA, TBB, Cilk, OpenACC, SYCL

I. INTRODUCTION

With the end of Dennard Scaling [1], multicore CPUs sharing a cache-coherent address space are ubiquitous. To exploit the parallelism available from multicore systems, programmers use multithreaded programming models. Programming models that support multithreaded execution include pThreads [2] for

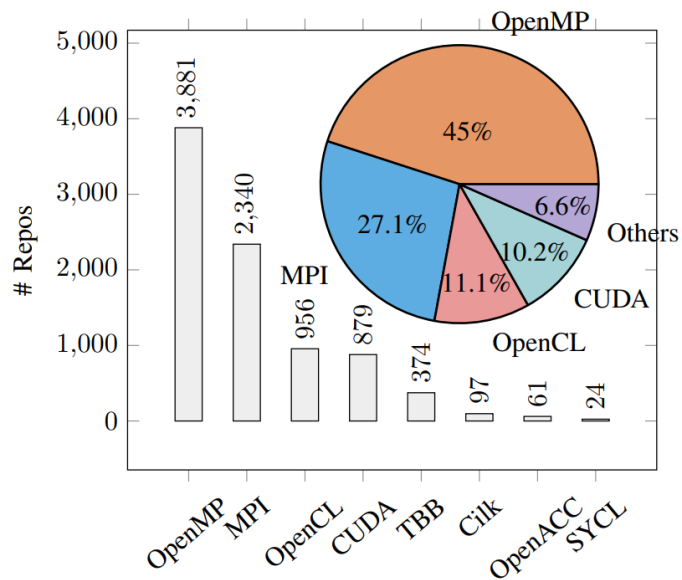
low-level and OS-level operations, TBB [3] or Cilk [4] for task-level parallelism in C++, and OpenMP [5] for directive-driven parallelism.

Despite the popularity of multithreaded models, little empirical data is available to assess the actual usage of the various programming constructs from these models. While anecdotal data and feedback from user-support teams at supercomputing centers exist [6], a large-scale analysis has yet to be conducted. In this paper, we perform a statistical analysis of repositories from GitHub to study the usage of parallel programming models. Our analysis reveals that OpenMP is the dominant programming model for writing multithreaded applications. We also go inside applications to gather usage data on specific OpenMP constructs. Finally, we consider the evolution of OpenMP and the adoption of newly included constructs as OpenMP Specifications are released.

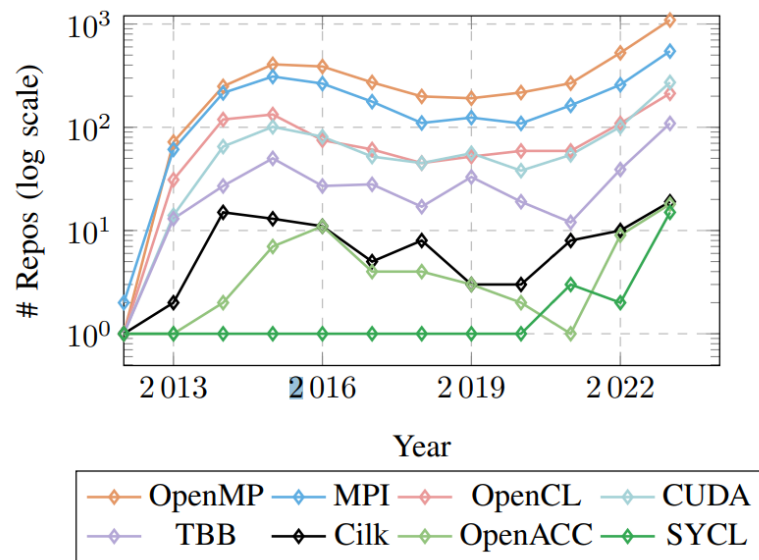
II. OPENMP FUNDAMENTALS VS. OTHER PARALLEL PROGRAMMING APIS

In this section, we provide a brief overview of OpenMP and other parallel programming APIs. OpenMP [7] defines a simple and portable approach to shared-memory parallel programming. OpenMP makes parallel programming more accessible by offering a directive-based approach, where directives inserted into the code guide the compiler as it generates parallel code. These directives provide high-level abstractions to specify parts of a code to execute in parallel. OpenMP uses the fork-join model of parallelism, where a single thread (the PRIMARY thread) on encountering a parallel directive forks a team of threads, each of which executes the code in a parallel region independently. Synchronization constructs, such as barriers, coordinate multithreaded execution, while shared variables facilitate data sharing among threads. OpenMP emphasizes portability with a standardized API that can be used across different platforms, hardware, and programming languages. These features make it easier to write parallel code that can be compiled and executed on systems supporting OpenMP.

HPCorpus: Parallel APIs

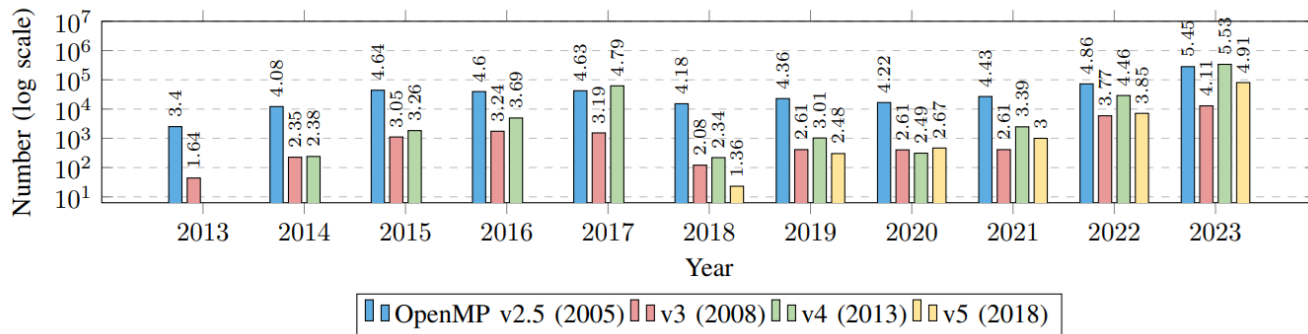
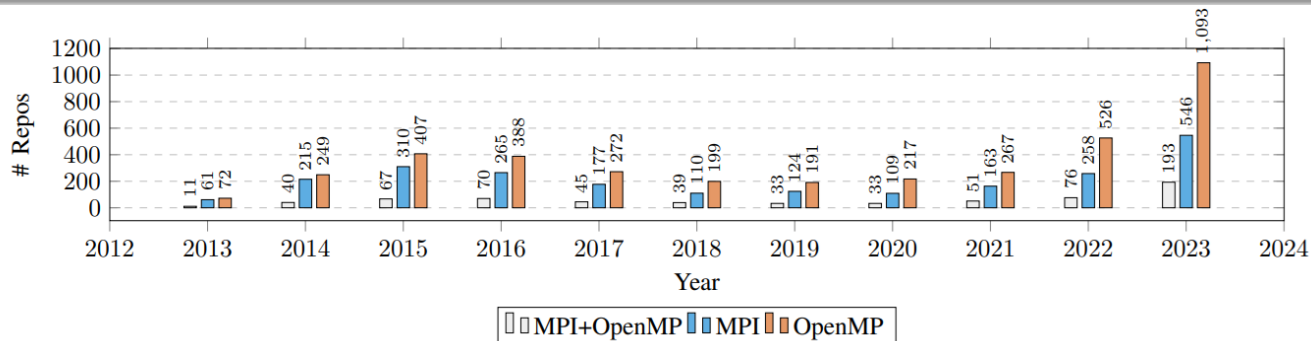


Parallel programming API usage in HPCorpus

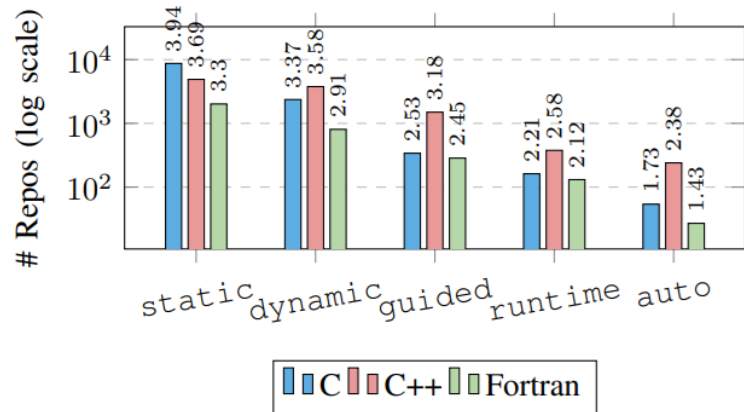
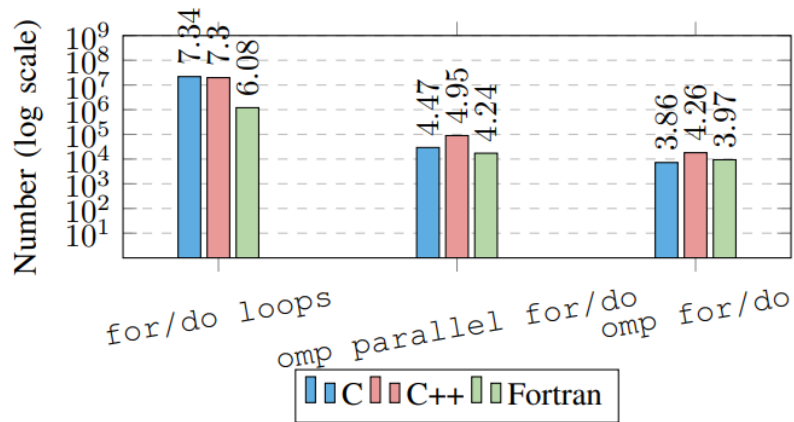


Usage trends in HPCorpus over the last decade

OpenMP and MPI usage breakdown



OpenMP loops and scheduling attributes

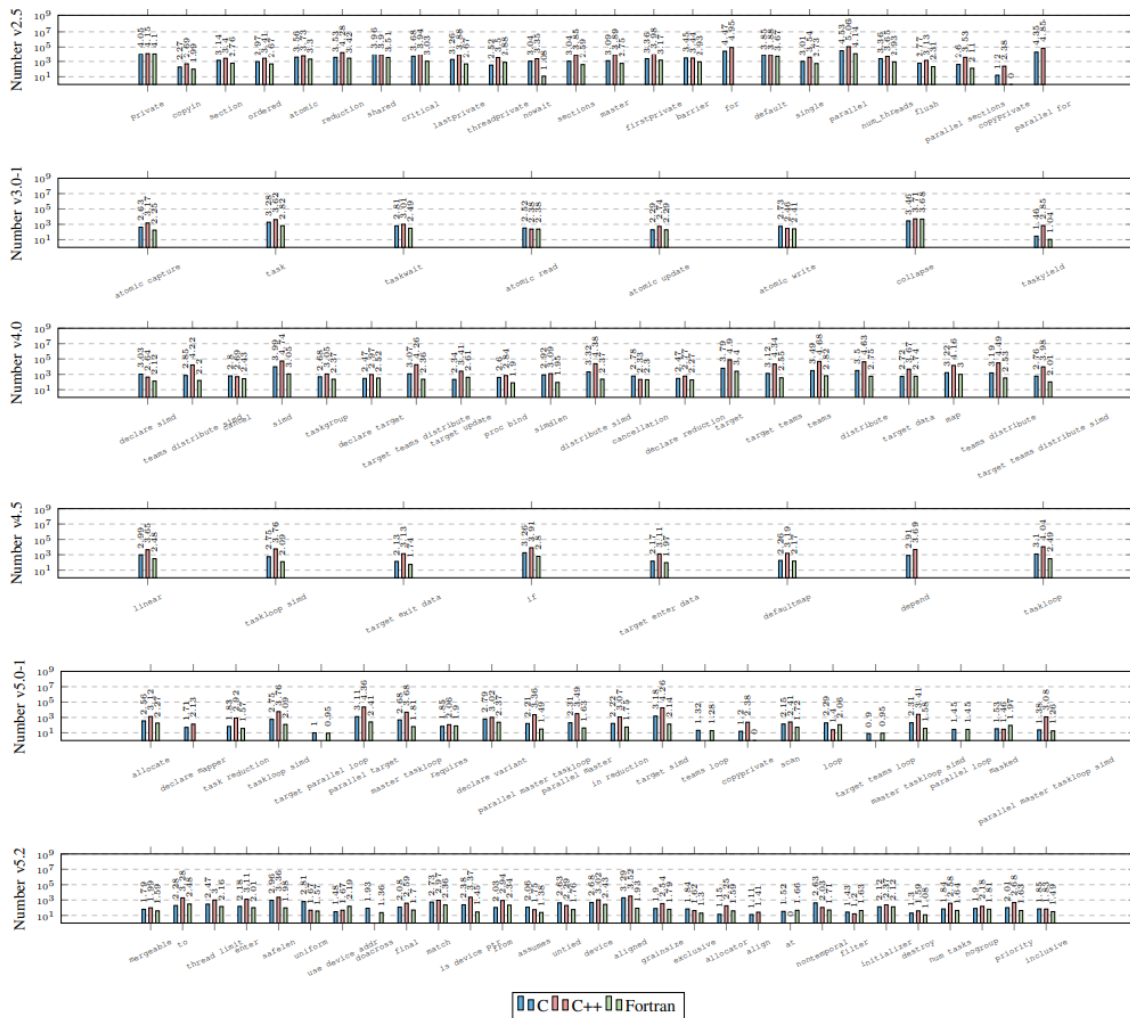


OpenMP meltdown

Each and every directive, clause etc. was carefully analysed for any future usage, including:

- SIMD
- teams distribute
- tasks
- Data-sharing attribute clauses

...



Scope is all you need:

Intro and objectives - LLMs for HPC: Gal Oren

Novel HPC code database - HPCorpus: Niranjan Hasabnis

Compiler-oriented tokenizer - Tokompiler: Vy A. Vo

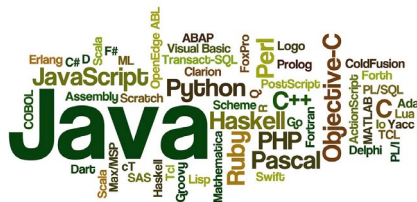
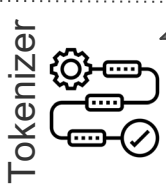
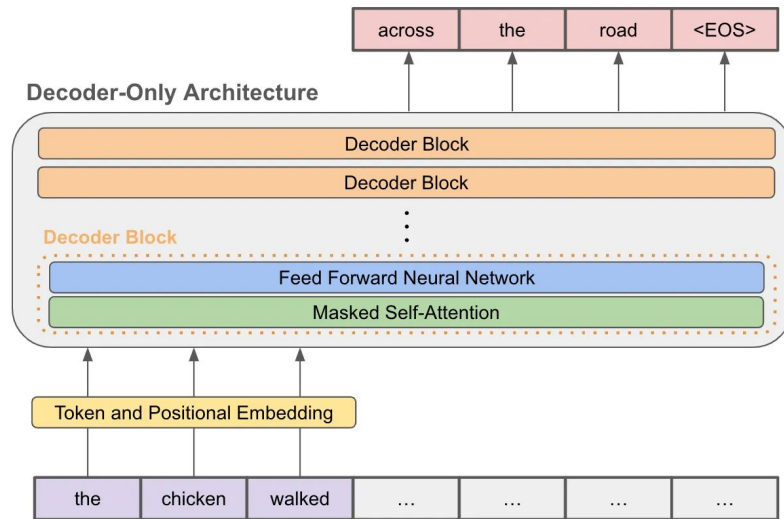
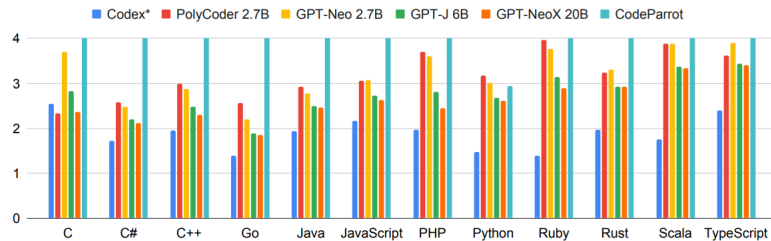
Downstream Task #1 - OpenMP: Tal Kadosh

Downstream Task #2 - MPI: Nadav Schneider

Pre-train or fine-tune?

- Why do we need an LLM trained on Java or Python (i.e., PolyCoder, StarCoder) for HPC-specific tasks?
- Can we use non-trivial or domain-specific tokenization for HPC tasks?
- Can domain-specificity improve the computational and therefore financial efficiency of pretraining?

Xu, Frank F., et al. "A systematic evaluation of large language models of code." *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*. 2022.



3

Perplexity

2

Pre-train/
Fine-tune

1

Corpora

Current Common Tokenization

- Current tokenization strategies for NL likely do not address the needs of HPC tasks
- Can we do better if we will find the way to:
 - Remove NL
 - Remove local semantics
 - Keep structure
 - Add more representations?

Tokens
336

Characters
508

GPT3 Tokenizer of the Pi C code:
<https://platform.openai.com/tokenizer>

```
#include <stdio.h>
int main() {
    int r[2800 + 1];
    int i, k;
    int b, d;
    int c = 0;
    for (i = 0; i < 2800; i++) {
        r[i] = 2000;
    }
    for (k = 2800; k > 0; k -= 14) {
        d = 0;
        i = k;
        for (;;) {
            d += r[i] * 10000;
            b = 2 * i - 1;
            r[i] = d % b;
            d /= b;
            i--;
            if (i == 0) break;
        }
    }
}
```

Tokens
206

Characters
508

Codex Tokenizer of the Pi C code:
<https://platform.openai.com/tokenizer>

```
#include <stdio.h>
int main() {
    int r[2800 + 1];
    int i, k;
    int b, d;
    int c = 0;
    for (i = 0; i < 2800; i++) {
        r[i] = 2000;
    }
    for (k = 2800; k > 0; k -= 14) {
        d = 0;
        i = k;
        for (;;) {
            d += r[i] * 10000;
            b = 2 * i - 1;
            r[i] = d % b;
            d /= b;
            i--;
            if (i == 0) break;
        }
    }
}
```

Tokompiler

- **Generate Replaced Code:**

Create a version of the original code with anonymized variable names, numbers, and strings

- **AST Parsing:**

Parse the anonymized code using tree-sitter or any suitable parser to generate an AST

- **Recreate AST Changes:**

Update the AST to reflect changes made during anonymization. Keep a dictionary of all changes done per file/function to facilitate restoring semantics back later

- **AST to Code-Tokenize:**

Transform the updated AST back into code, eliminating any comments, new lines, and READMEs that may have been introduced during anonymization. This code-tokenized version will have a much smaller number of tokens

- **Token Splitting:**

Split multi-part tokens (e.g., "var_1" to ["var", "1"]) to ensure that the model comprehends variable names as a combination of type and a unique identifier

- **Random Number Attachment:**

For recurrent tokens (e.g., "var_1" or "num_2"), use statistics to attach random numbers from a predefined range (e.g., 1 to 1000) during each tokenization. The attached numbers are randomly chosen without any relation to the type or order of the replaced tokens or the file/function length

HPCorpus-Fortran Tokompiler

Tokenized code

```
[['function', 'func', '48', '(', 'arg', '128', ',', 'arg', '807', ')', 'result',  
'(', 'func', '180', ')', 'implicit', 'none', 'integer', ',', 'intent', '(', 'in',  
)', ',', 'arg', '128', ',', 'arg', '807', 'real', '(', 'num', '5', ')', ',',  
'func', '180', 'real', '(', 'num', '5', ')', ',', 'var', '377', ',', 'var',  
'84', ',', 'var', '967', 'integer', ',', 'var', '821', 'external', ',',  
'func', '123', 'integer', ',', 'var', '63', 'var', '63', '=', 'num', '156',  
'call', 'var', '719', '(', 'arg', '807', ')', 'do', 'var', '821', '=', 'num',  
'315', ',', 'arg', '128', 'var', '84', '=', 'func', '123', '(', ')', '*',  
'num', '357', ',', 'num', '315', 'var', '967', '=', 'func', '123', '(', ')',  
'*', 'num', '357', ',', 'num', '315', '(', 'var', '84', '*', 'var', '84',  
'+', 'var', '967', '*', 'var', '967', '<', 'num', '315', ')', 'then', 'var',  
'63', '=', 'var', '63', '+', 'num', '315', 'end', 'var', '377', '=',  
'num', '539', '*', 'func', '937', '(', 'var', '63', ')', '/', 'func', '937',  
'(', 'var', '821', ')', 'end', 'do', 'func', '180', '=', 'num', '539', '*',  
'func', '937', '(', 'var', '63', ')', '/', 'func', '937', '(', 'arg', '128', ')',  
'end', 'function']
```

687.72

Tokenized AST

```
'translation_unit', '(', 'declaration', '(', 'function',  
'function_declarator', '(', 'func', '48', 'parameter_list', '(', '(',  
'parameter_declaration', '(', 'arg', '128', ')', ',',  
'parameter_declaration', '(', 'arg', '807', ')', ')', ')', ')',  
'declaration', '(', 'macro_type_specifier', '(', 'result', '(', 'type',  
'descriptor', '(', 'func', '180', ')', ')', ')', 'implicit', ')', 'declaration',  
'(', 'none', 'integer', ',', 'function_declarator', '(', 'intent',  
'parameter_list', '(', '(', 'parameter_declaration', '(', 'in', ')', ')',  
)', ')', '(', ',', 'arg', '128', ')', ',', '(', 'arg', '807',  
'function_declarator', '(', 'real', 'parameter_list', '(', '(',  
'parameter_declaration', '(', 'num', '5', ')', ')', ')', ',', ',',  
'func', '180', 'function_declarator', '(', 'real', 'parameter_list',  
'(', '(', 'parameter_declaration', '(', 'num', '5', ')', ')', ')', ')', ',', ',',  
'var', '377', ',', 'var', '84', ',', '(', 'var', '967', 'integer', ',',  
'var', '821', 'external', ',', ',', 'func', '123', 'integer', ',', ',', 'var',  
'63', ')', 'init_declarator', '(', 'var', '63', '=', 'num', '156', ')', ')',  
'declaration', '(', 'call', 'function_declarator', '(', 'var', '719',  
'parameter_list', '(', '(',  
...
```

1099.1

177 primitives + 1000 locals = 1177

Pre-training on HPCorpus Fortran

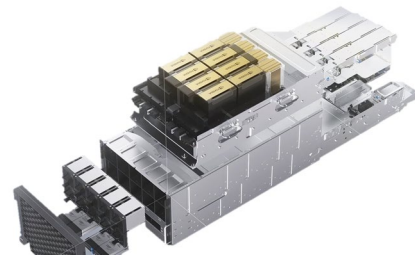
- New given form of data (stripped of NL and only using structured blocks of code) help the models in improving the results
- The usage of Tokompiler increased performance by 230%
- Training speed improved by an average of 10% (as vocabulary and tokenized corpus size are smaller, 1177 against 50K).

HPCorpus Fortran Pre-training

	PolyCoder		SPT-Code	
	BPE	Tokompiler	NLTK	Tokompiler
Code	2.46	1.60	2.50	1.08
Code+AST	NA	NA	2.72	1.08
Model Size	162M	88M	238M	201M
Time to train (min)	435	390	186	150



4 x A40 48GB GPUs



8 x H100 80G GPUs

Pre-training on HPCorpus C/C++

- We pretrained an LLM with a standard size with both tokenizers
- On architectures that only vary in the size of the embedding, performance is better with Tokompiler
- Reducing the size of the Tokompiler model yields better performance

HPCorpus C/C++ Pre-training

	PolyCoder		
	BPE	Tokompiler	Tokompiler-small
Code	1.93	1.71	1.65
Code+AST	NA	NA	NA
Model Size	2.8B	2.5B	638M
Time to train (min)	8300	8125	6386



4 x A40 48GB GPUs

Scope is all you need:

Intro and objectives - LLMs for HPC: Gal Oren

Novel HPC code database - HPCorpus: Niranjan Hasabnis

Compiler-oriented tokenizer - Tokompiler: Vy A. Vo

Downstream Task #1 - OpenMP: Tal Kadosh

Downstream Task #2 - MPI: Nadav Schneider

Our Vision

Does a language model can answer the following questions?

- Can we detect opportunity for parallelization?
- What should the variables scoping be?
- Can we completely generate a pragma given an input?

```
void dotProduct(int* A, int* B, int size) {  
    int sum = 0;  
    for (int i = 0; i < size; ++i) {  
        sum += A[i] * B[i];  
    }  
    return sum;  
}
```



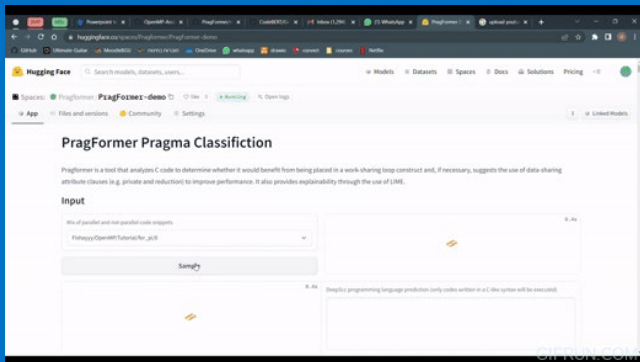
Large Language Model



```
void dotProduct(int* A, int* B, int size) {  
    int sum = 0;  
    #pragma omp target data map(to: A[0:size], B[0:size]) map(from: sum){  
        #pragma omp target teams distribute parallel for reduction(+:sum)  
        for (int i = 0; i < size; ++i) {  
            sum += A[i] * B[i];  
        }  
    }  
    return sum;  
}
```

PragFormer

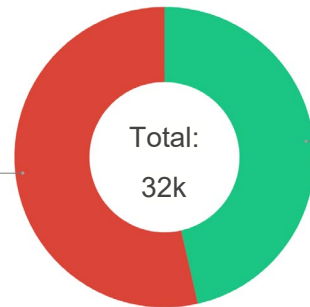
- Can we detect opportunity for parallelization?
- Is it possible to classify the need for OpenMP pragma?



HP Corpus



w/o OpenMP
53.6%



with OpenMP
46.4%

```
int print_dotProd(double *A, double *B) {
```

```
    for (int i = 0; i < size; ++i) {  
        printf("%f * %f", A[i], B[i]);  
    }
```

Negative Example

```
#pragma omp reduction(+:sum)
```

```
for (int i = 0; i < size; ++i) {  
    sum += A[i] * B[i];  
}
```

Positive Example

```
}
```


PragFormer

- Can we detect opportunity for parallelization?
- Is it possible to classify the need for OpenMP pragma?

```
for (int i = 0; i < size; ++i) {  
    sum += A[i] * B[i];  
}
```



PragFormer
Pragma

PragFormer
Private

PragFormer
Reduction



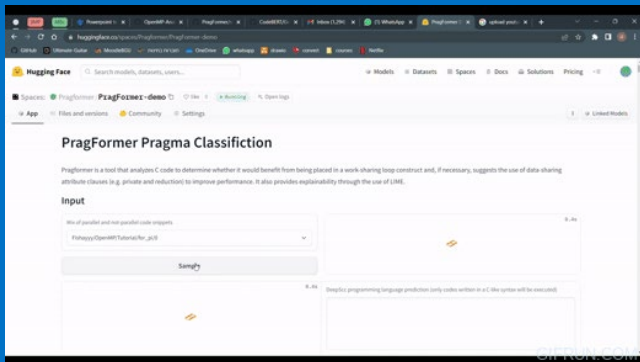
Yes / No



Yes / No

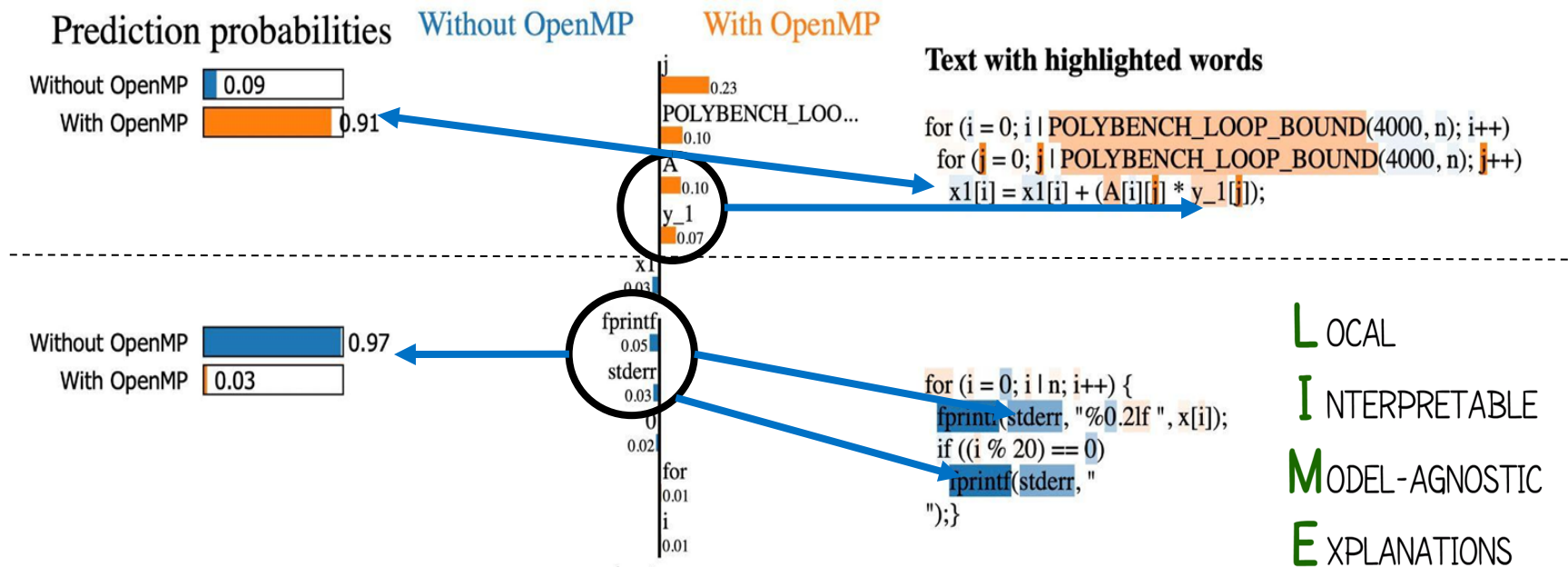


Yes / No



PragFormer

Strongly influenced by local semantics



Can We Do Better?

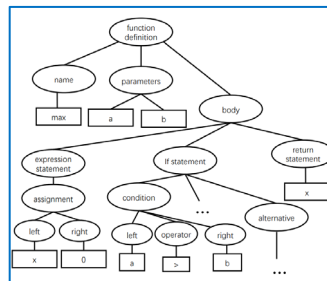
- Exploit different code representations
- Provide extended scope of code
- Utilize an LLM pre-trained on codes
- Use a unified model

Code

```
def max(a, b):  
    x=0  
    if b>a:  
        x=b  
    else:  
        x=a  
    return x
```

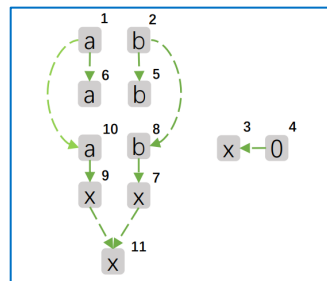
Syntactic
Information

AST



Structural
Information

DFG

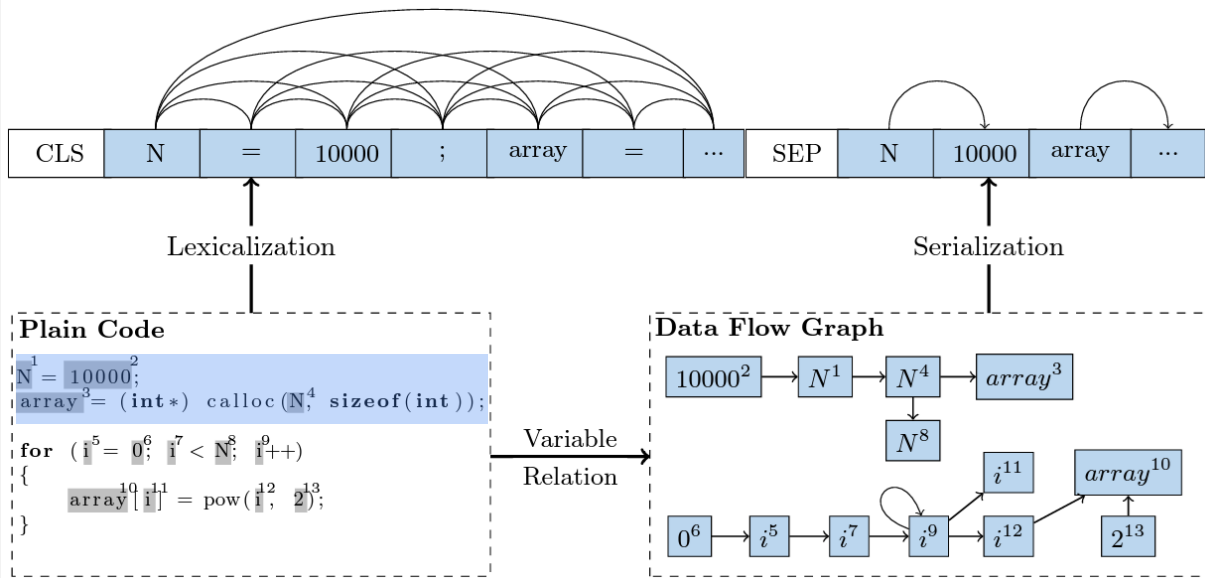


Semantic
Information

OMPify

Input Format

- Ompify classifies the need for the parallel for pragma and its clauses simultaneously
- It is utilizing the extended scope of the for loop as supplementary information
- It is using DFG as an additional code representation



OMPify

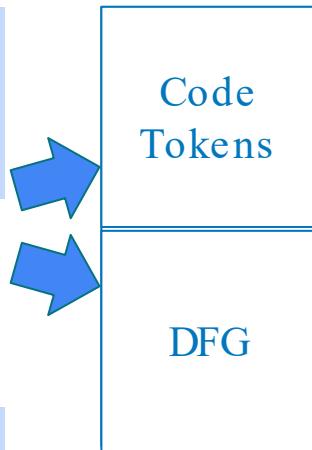
Unified model for multiple code representations

Multiple Code Representations

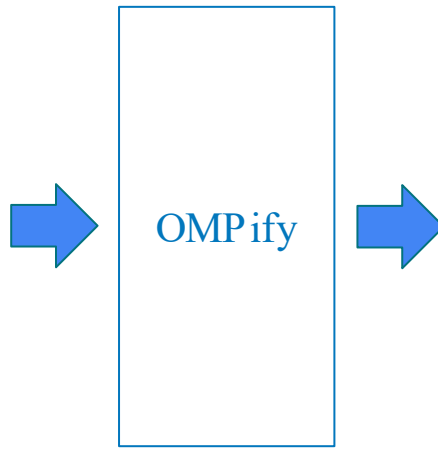
```
Int size = 10000;  
int* A = (int*)malloc(sizeof(int) * size);  
int* B = (int*)malloc(sizeof(int) * size);  
int sum = 0;
```

```
for (int i = 0; i < size; ++i) {  
    sum += A[i] * B[i];  
}
```

```
printf("sum = %d", sum);
```



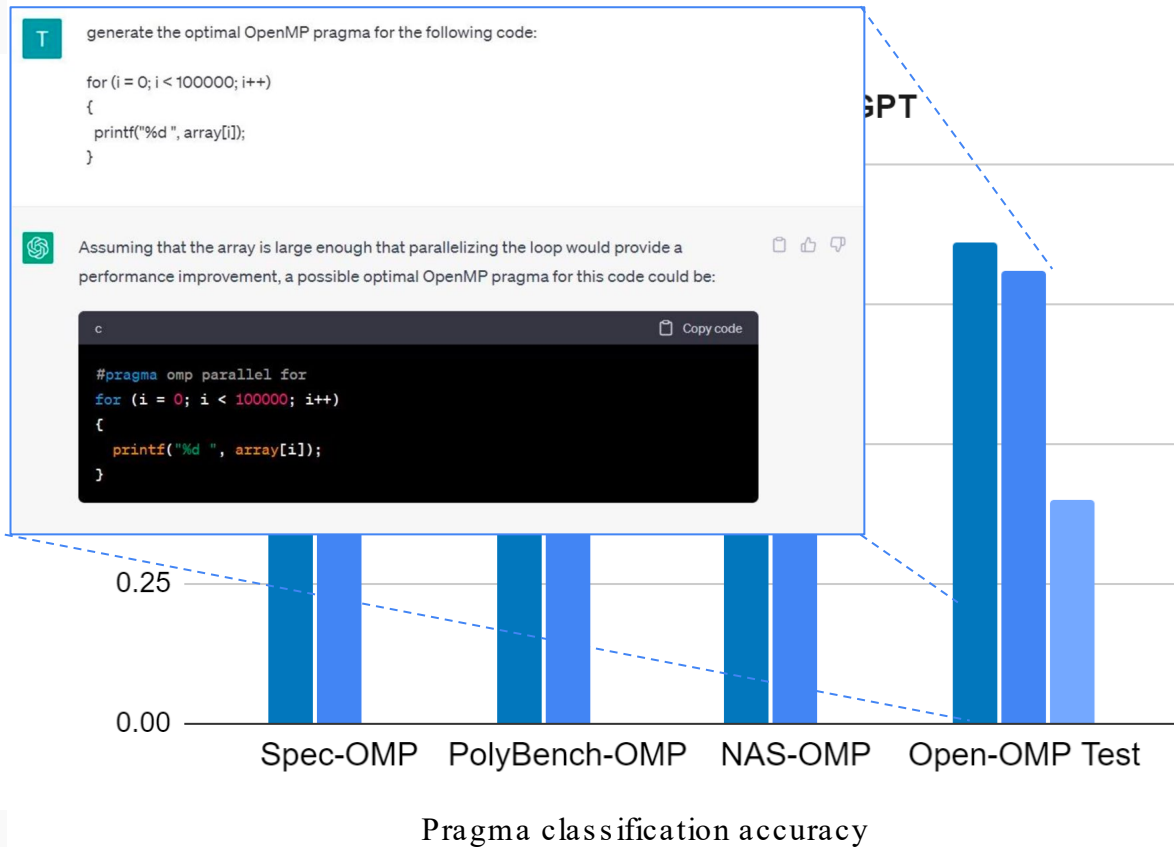
Unified Model



Yes / No Pragma
Yes / No Private
Yes / No Reduction

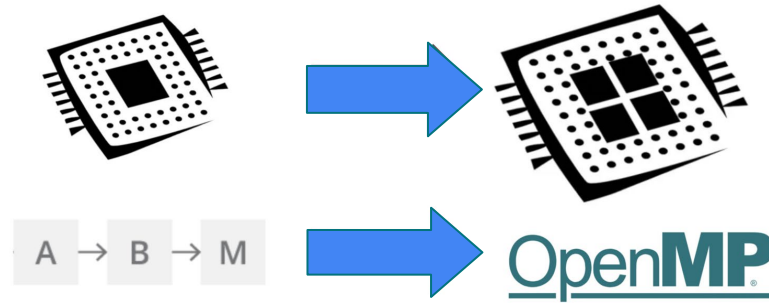
Results

- The data flow graph is found to be the most influential representation.
- The extended scope provides a slight improvement.
- Using a unified model is found to be superior to using separate models.



What About Pragma Generation?

```
for (int i = 0; i < size; ++i) {  
    sum+= A[i] * B[i];  
}
```



```
#pragma omp parallel for  
reduction(+:sum)
```



```
for (int var_64 = num_5; var_64 < var_826; ++var_64) {  
    var_3+= arr_26[var_64] * arr_109[var_64];  
}
```



```
#pragma omp parallel for  
reduction(+:var_3)
```

Scope is all you need:

Intro and objectives - LLMs for HPC: Gal Oren

Novel HPC code database - HPCorpus: Niranjan Hasabnis

Compiler-oriented tokenizer - Tokompiler: Vy A. Vo

Downstream Task #1 - OpenMP: Tal Kadosh

Downstream Task #2 - MPI: Nadav Schneider

Distributed Parallelism-Domain Decomposition

CHALLENGES

- Is parallelization possible?
- Comprehensive code understanding
- Misplace send/receive functions
- Selecting the right functions in the right locations

MPI programming poses challenges!

```
int min(arge, argv)
{
    int done = 0, n = 10000, rank, numprocs, i;
    double pi_total, pi, h, sum_x;

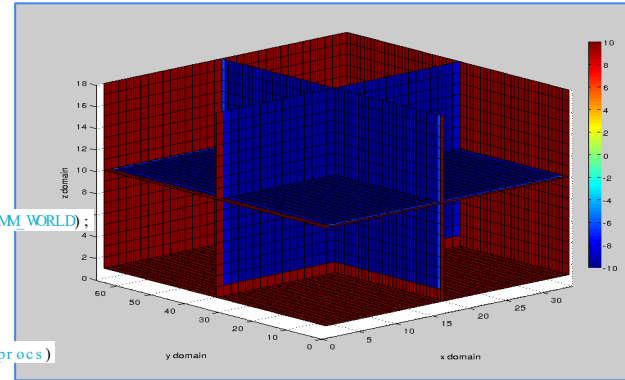
    MPI_Init(&arge, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    while (!done)
    {
        MPI_Bcast(&n, 1, MPI_INT, 0, MPI_COMM_WORLD);
        if (n == 0) break;

        h = 1.0 / (double) n;
        sum = 0.0;

        for (i = rank + 1; i <= n; i += numprocs)
        {
            x = h * ((double) i - 0.5);
            sum += 4.0 / (1.0 + x*x);
        }

        pi_total = h * sum;
        MPI_Reduce(&pi_total, &pi, 1, MPI_DOUBLE, MPI_SUM 0,
MPI_COMM_WORLD);
    }
    MPI_Finalize();
    return 0;
}
```



Serial vs. Parallel Code

OpenMP
Pragma based

```
#pragma omp parallel for private(i) reduction(+:sum)  
for (int i = 0; i < size; ++i) {  
    sum += A[i] * B[i];  
}
```

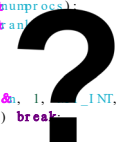
Serial Code

Parallel Code

MPI
Function based

```
int main(argc, argv)  
{  
    int done = 0, n = 10000, rank, numprocs, i;  
    double pi_total, pi, h, sum, x;  
  
    MPI_Init(&argc, &argv);  
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);  
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);  
  
    while (!done)  
    {  
        MPI_Bcast(&n, 1, MPI_INT, 0, MPI_COMM_WORLD);  
        if (n == 0) break;  
  
        h = 1.0 / (double) n;  
        sum = 0.0;  
  
        for (i = rank; i < n; i += numprocs)  
        {  
            x = h * ((double) i - 0.5);  
            sum += 4.0 / (1.0 + x*x);  
        }  
  
        pi_total = h * sum;  
        MPI_Reduce(&pi_total, &pi, 1, MPI_DOUBLE, MPI_SUM, 0,  
MPI_COMM_WORLD);  
    }  
    MPI_Finalize();  
    return 0;  
}
```

Serial Code



MPI ROADMAP (For Automatic Parallelization)

Q1

SCOPE

What is the model inputs?
Repositories?
For loops?
MPI_INIT to
MPI_FINALIZE
constructs?

Q2

SERIALIZATION

How to serialize the
code: MPI Variables
removal?
MPI Functions
removal?
Rule based model?

Q3

TASK

Do we train a
Generative or a
Classification
model

MPI-rical - Partial Solution

A1

SCOPE

What is the model
inputs?
MPI_INIT to
MPI_FINALIZE
constructs?

A2

SERIALIZATION

How to serialize the
code:
MPI Functions removal

A3

TASK

Train a
Generative
model

MPI-rical Partial Solution

A1

SCOPE

What is the model
inputs?
MPI_INIT to
MPI_FINALIZE
constructs?

```
int min(argc, argv)
{
    int done = 0, n = 10000, rank, numprocs, i;
    double pi_total, pi, h, sum, x;

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    while (!done)
    {
        MPI_Bcast(&n, 1, MPI_INT, 0, MPI_COMM_WORLD);
        if (n == 0) break;

        h = 1.0 / (double) n;
        sum = 0.0;

        for (i = rank + 1; i <= n; i += numprocs)
        {
            x = h * ((double) i - 0.5);
            sum += 4.0 / (1.0 + x*x);
        }

        pi_total = h * sum;
        MPI_Reduce(&pi_total, &pi, 1, MPI_DOUBLE, MPI_SUM, 0,
MPI_COMM_WORLD);
    }
    MPI_Finalize();
    return 0;
}
```

MPI-rical Partial Solution

A2

SERIALIZATION

How to serialize the
code:
MPI Functions removal

```
int min(argc, argv)
{
    int done = 0, n = 10000, rank, numprocs, i;
    double pi_total, pi, h, sum, x;

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    while (!done)
    {
        MPI_Bcast(&n, 1, MPI_INT, 0, MPI_COMM_WORLD);
        if (n == 0) break;

        h = 1.0 / (double) n;
        sum = 0.0;

        for (i = rank + 1; i <= n; i += numprocs)
        {
            x = h * ((double) i - 0.5);
            sum += 4.0 / (1.0 + x*x);
        }

        pi_total = h * sum;
        MPI_Reduce(&pi_total, &pi, 1, MPI_DOUBLE, MPI_SUM, 0,
MPI_COMM_WORLD);
    }
    MPI_Finalize();
    return 0;
}
```

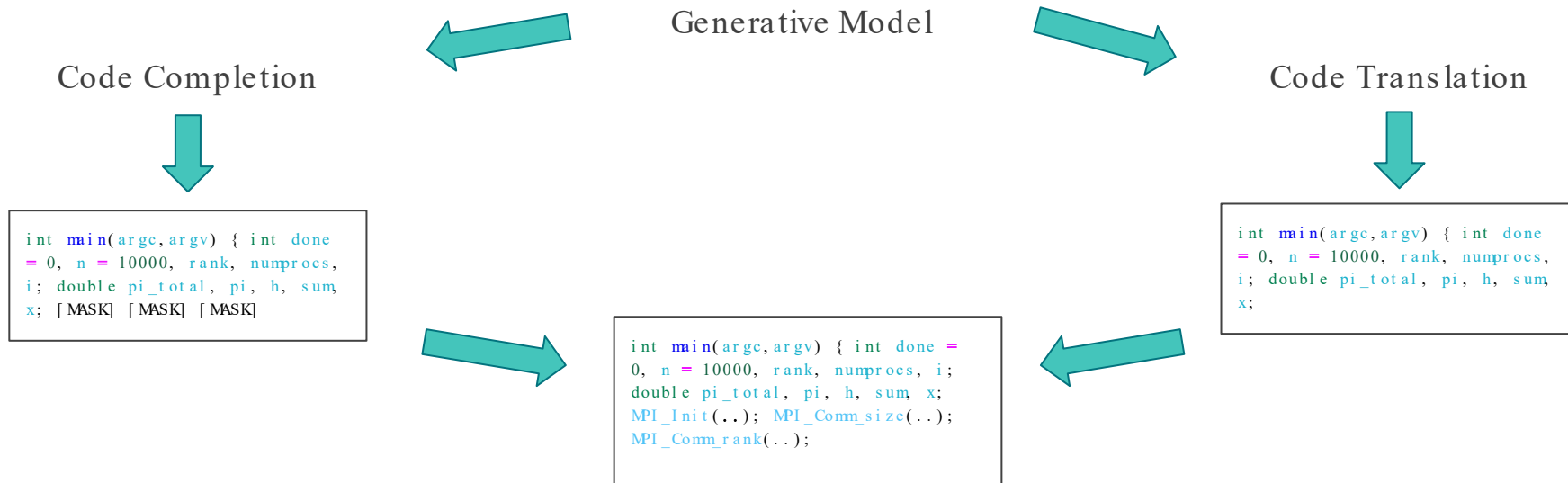
SEMI serial code

A3

TASK

Train a
Generative
model

MPI-rical - Partial Solution



MPI-rical Partial Solution

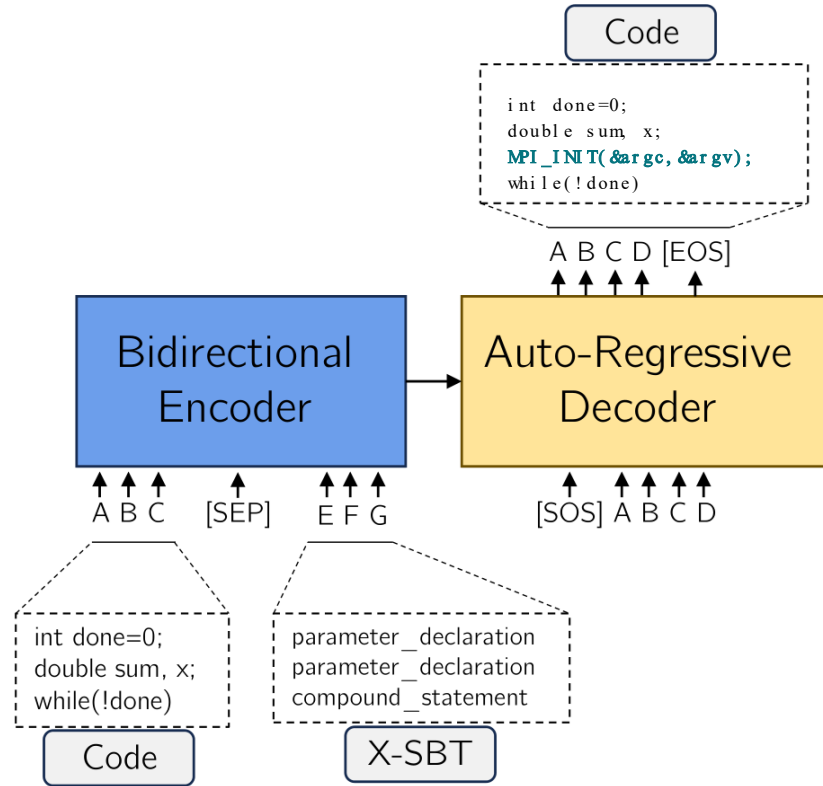
A3

TASK

Train a
Generative
model

Code Translation

MPI parallel code



SEMI serial code

Dataset

- SCOPE: Init to Finalize Constructs.
- 25k SPMD and 25k non-SPMD.
- NEGATIVE EXAMPLES: created by MPI functions removal.



```
int min(argc, argv)
{
    int done = 0, n = 10000, rank, numprocs, i;
    double pi_total, pi, h, sum, x;

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    while (!done)
    {
        MPI_Bcast(&n, 1, MPI_INT, 0, MPI_COMM_WORLD);
        if (n == 0) break;

        h = 1.0 / (double) n;
        sum = 0.0;
    }
}
```

Positive Example

```
int min(argc, argv)
{
    int done = 0, n = 10000, rank, numprocs, i;
    double pi_total, pi, h, sum, x;

    while (!done)
    {
        if (n == 0) break;

        h = 1.0 / (double) n;
        sum = 0.0;
    }
}
```

Negative Example

Dataset

- SCOPE: Init to Finalize Constructs.
- 25k SPMD and 25k non-SPMD.
- NEGATIVE EXAMPLES: created by MPI functions removal.

# Line	Amount
≤ 10	2,670
11-50	22,361
51-99	14,078
≥ 100	10,575

Function	Amount
MPI_Finalize	35,983
MPI_Comm_rank	32,312
MPI_Comm_size	28,742
MPI_Init	25,114
MPI_Recv	10,340
MPI_Send	9,841
MPI_Reduce	8,503
MPI_Bcast	5,296



**F,
C,C++
files**



**MPI_INIT to
MPI_FINALIZE**

MPI-rical - Partial Solution

```
int min(argc, argv)
{
    int done = 0, n = 10000, rank, numprocs, i;
    double pi_total, pi, h, sum, x;

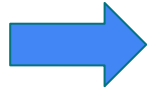
    while (!done)
    {
        if (n == 0) break;

        h = 1.0 / (double) n;
        sum = 0.0;

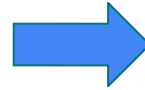
        for (i = rank + 1; i <= n; i += numprocs)
        {
            x = h * ((double) i - 0.5);
            sum += 4.0 / (1.0 + x*x);
        }

        pi_total = h * sum;
    }

    return 0;
}
```



MPI-rical
Can it be
SPMD
parallelized?



MPI Function Classification
(MPI_Init(...), MPI_Comm_rank(...), MPI_Send(...))

MPI Function Location



Yes / No

Results

Generation as Classification

Ground-Truth: `int main(argc,argv){ int done = 0, rank, numprocs; MPI_Init(..); MPI_Comm_size(..); MPI_Comm_rank(..); while(!done){`
Prediction: `int main(argc,argv){ int done = 0, rank, numprocs; MPI_Init(..); MPI_Bcast(..); while(!done){`

TP FP FN

Code	M-F1	M-Precision	M-Recall
Array Average	0.88	1.0	0.8
Vector Dot Product	0.88	1.0	0.8
Min-Max	0.66	1	0.5
Matrix-Vector Multiplication	0.9	0.83	1.0
Sum (Reduce & Gather)	0.8	1.0	0.6
Merge Sort	1.0	1.0	1.0
Pi Monte-Carlo	1.0	1.0	1.0
Pi Riemann Sum	1.0	1.0	1.0
Factorial	0.88	1.0	0.8
Fibonacci	1.0	1.0	1.0
Trapezoidal Rule (Integration)	1.0	1.0	1.0
Total	0.91	0.98	0.86

Quality Measure	<i>MPICodeCorpus</i>
M-F1	0.87
M-Precision	0.85
M-Recall	0.89
MCC-F1	0.89
MCC-Precision	0.91
MCC-Recall	0.87
BLEU	0.93
Meteor	0.62
Rouge-l	0.95
ACC	0.57

What Next?

Full MPI Program Generation

```
int min(argc,argv)
{
  int done = 0, n = 10000, rank, numprocs, i;
  double pi_total, pi, h, sum, x;

  while (!done)
  {
    if (n == 0) break;
    h = 1.0 / (double) n;
    sum = 0.0;
```

Tokompiler

```
int min(argc,argv)
{
  int var_1 = num_0, var_2 = num_1, var_2, var_3, var_4;
  double var_5, var_6, var_7, var_8, var_9;

  while (!var_1)
  {
    if (var_2 == num_0) break;
    var_7 = num_2 / (double) var_2;
    var_8 = num_3;
```

MPI Function Classification

(MPI_Init(&argc, &argv), MPI_Comm_rank(MPI_COMM_WORLD, &rank) ...)

```
int min(argc,argv)
{
  int done = 0, n = 10000, rank, numprocs, i;
  double pi_total, pi, h, sum, x;

  MPI_Init(&argc, &argv);
  MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
  MPI_Comm_rank(MPI_COMM_WORLD, &rank);

  while (!done)
  {
    MPI_Bcast(&done, 1, MPI_INT, 0, MPI_COMM_WORLD);
    if (n == 0) break;
    h = 1.0 / (double) n;
    sum = 0.0;
```

Tokompiler

```
int min(argc,argv)
{
  int var_1 = num_0, var_2 = num_1, var_2, var_3, var_4;
  double var_5, var_6, var_7, var_8, var_9;
  func_339(@type_7, @type_5);
  func_459(var_23, @var_18);
  func_113(var_43, @var_11);

  while (!var_1)
  {
    func_287(@var_29, num_43, var_169, num_0, var_95);

    if (var_2 == num_0) break;

    var_7 = num_2 / (double)
var_2;
    var_8 = num_3;
```

Comp-Coder

Scope is all you need:

Intro and objectives - LLMs for HPC: Gal Oren

Novel HPC code database - HPCorpus: Niranjan Hasabnis

Compiler-oriented tokenizer - Tokompiler: Vy A. Vo

Downstream Task #1 - OpenMP: Tal Kadosh

Downstream Task #2 - MPI: Nadav Schneider

Thanks!

Contact us:

galoren@cs.technion.ac.il

```
2, 19.3m3EM *  
y, {align}m?2  
oms8? ;{f}a  
on, l>msa  
sbno03aa  
na ;{ba  
one bol  
em l>  
e + "  
qm  
t
```

```
or(int i = 0; i < Side  
m + 1; i++)render.GetComponent  
.. LineRenderer>().SetPosition(0, ListPo  
float angle = i * 2 * Mathf.PI / SideNum; using  
nityEngine;for ListPos.Add(new Vector3(Radius *  
athf.Cos(angle), Radius * Mathf.Sin(angle), 0)); for (i  
t i = 0; i < ListPos.Count - 1; i++)render.GetComponent  
<LineRenderer>().SetPosition(0, ListPos[i]); GameObj  
ect line = render <LineRenderer>().SetPosit ion(0, List  
P a[i]);}() public// Debug.Log( Renderer>().SetPositio  
n(0, ListPos[start]); fo for (int i = 0; i ++)  
mponent for (int i = 0; i < Side Num + 1; i++)re nder.GetCent  
LineRener/Collections;Setart = + start); using UnityEng  
ine;for // DebUse this for initialization void Start() Start(  
onoBehaviour {SetColor("int late is call float angle  
= i athf. GameObject line = GameObject line = In  
stantit(Li ePrefab);puc List <Vectr3> List os; L  
ined); // Use t id Start() render.GetComponR  
derer render = line. mponent<Linentiate(L  
in ePref
```

```
tiat  
istPos;  
render <Line  
Renderer>().S  
for (int i = 0; i < S  
start = " + start); us  
Start() MonoBehavio  
GameObject line = Gar  
= " + en d); // Use t id S  
ePref neRenderer rende  
s[]];= (start + 3) % SideN  
anti int end = (start + 3) %  
// Debug.Log("start = "  
er.GetComponent<LineRe  
d eN m i line.Get  
for (i nt i = 0; i < tP flo  
g("end = " + end); ren  
.SetPositi on(0, L istPo  
; GameObject line = Insta  
der.GetComponent<LineRen  
orstPos.Count - 1; i++)rend for  
void Upd ate() GameObject line  
LineRenderer>()); int start = i; int e
```