

Compilation and Backend-Independent Vectorization for Multi-Party Computation (Extended Version)

Benjamin Levy and Ben Sherman*
{levyb3,shermb}@rpi.edu
Rensselaer Polytechnic Institute
Troy, New York

Muhammad Ishaq
ishaqm@purdue.edu
Purdue University
West Lafayette, Indiana

Lindsey Kennard
fireelemental.ne@gmail.com
STR
Boston, Massachusetts

Ana Milanova
milanova@cs.rpi.edu
Rensselaer Polytechnic Institute
Troy, New York

Vassilis Zikas
vzikas@purdue.edu
Purdue University
West Lafayette, Indiana

ABSTRACT

Research on MPC programming technology largely falls at the two ends of the classical compiler: (1) work on front-end language design (e.g., Wysteria, Viaduct) and (2) work on back-end protocol implementation (e.g., ABY, MOTION).

In this work, we formalize the MPC Source intermediate language and advance the idea of what we call backend-independent optimizations, in a close analogy to machine-independent optimizations in the classical compiler. We present a compiler framework that takes a Python-like routine and produces MOTION code. We focus on a specific backend-independent optimization: novel SIMD-vectorization on MPC Source, which we show leads to significant speedup in circuit generation time and running time, as well as significant reduction in communication and number of gates over the corresponding iterative schedule.

1 INTRODUCTION

Multi-party computation (MPC) allows N parties p_1, \dots, p_N to perform a computation on their private inputs securely. Informally, security means that the secure computation protocol computes the correct output (correctness) and it does not leak any information about the individual party inputs, other than what can be deduced from the output (privacy).

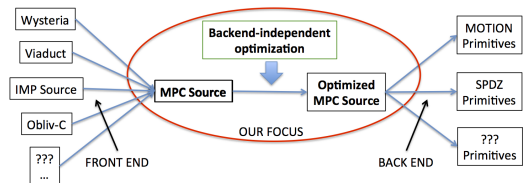
MPC theory dates back to the early 1980-ies [39, 23, 7, 15]. Long the realm of theoretical cryptography, MPC has seen significant advances in programming technology in recent years. These advances bring MPC closer to practice and wider applicability — MPC technology has been employed in real-world scenarios such as auctions [11], biometric identification [9], and privacy-preserving machine learning [31, 30]. The goal is to improve technology so that programmers can write *secure* and *efficient* programs without commanding extensive knowledge of cryptographic primitives.

The problem, therefore, is to build a high-level programming language and a compiler, and there has been significant advance in this space, e.g., [6, 10, 41, 24, 28, 1, 12] among other work. Current research largely falls at the two ends of the classical compiler: (1) work on *front-end* language design and (2) work on *back-end* protocol implementation. Work on language design focuses on

high-level constructs necessary to express multiple parties, computation by different parties, and information flow from one party to another [34, 1]. On the other end, work on protocol implementation focuses on cryptographic foundations and their efficient circuit-level implementation [18, 5, 12], e.g., implementation of operations (e.g., MUL, ADD) using different sharing protocols (Boolean or Arithmetic GMW [23] or Yao’s garbled circuits [39]), as well as efficient share conversion from one representation to another.

Earlier compilers did both back-end and front-end translation as their aim was to demonstrate applicability of MPC on real-world programming problems. As the field advanced, works have focused more closely on front-end language design (e.g., Wysteria [34] and Viaduct [1]) or back-end “circuit-level” design and implementation (e.g., MOTION [12]).

In this work we focus on an intermediate language and what we call *backend-independent optimizations*, in a close analogy to *machine-independent optimizations* in the classical compiler. The following figure summarizes our key idea:



We formalize the MPC Source [25] intermediate representation and emphasize optimization over MPC Source. As in classical compilers, we envision different front ends (e.g., our front end IMP Source, other) compiling into MPC Source. MPC Source is particularly suitable for optimizations such as protocol mixing [14, 25, 21] and SIMD-vectorization, which takes advantage of amortization at the circuit level. The MPC Source IR exposes the *linear structure* of MPC programs, which simplifies program analysis; this is in contrast to source, which has if-then-else constructs. In the same time, MPC Source is sufficiently “high-level” to support analysis and optimizations that take into account control and data flow in a specific program. MPC Source facilitates cost modeling. Again as in classical compilers, we envision translation of MPC Source (optimized or unoptimized) into MOTION, SPDZ, or other back-end code.

*Joint first authors.

1.1 Our Contribution

In this paper, we develop a compiler framework that takes a Python-like routine and produces MOTION code: we describe (a) the IMP Source language, its syntax and semantic restrictions, (b) translation into MPC Source, (c) a specific backend-independent optimization: novel SIMD-vectorization on MPC Source, and (d) translation from MPC Source into MOTION code.

We focus on the MOTION framework as our back-end for several reasons. First, it is the state of the art in terms of performance [12]. Second, it provides an API over efficient implementation for a wide variety of cryptographic operations in three different protocols – Arithmetic GMW, Boolean GMW, and BMR – which allows for protocol mixing [14, 25, 21], a known backend-independent optimization. Third, MOTION provides API for SIMD-level operations, which amortize cost and lead to significant improvement in memory footprint and throughput [18, 5, 12]. It enables MPC Source-level vectorization, a key focus of this paper.

Our second contribution is an analytical model for cost estimation of amortized schedules. Originally, we hoped that optimal scheduling (under our model, which essentially minimizes the length of the schedule) was tractable, as the problem appeared simpler than the classical scheduling problem. Unfortunately, we show that optimal scheduling is NP-hard via a reduction to the Shortest Common Supersequence (SCS) problem. Cost modeling is important as it drives not only vectorization but optimizations such as protocol mixing and scheduling as well [14].

Our most important contribution is the implementation and evaluation of the compiler framework. We demonstrate expressivity of the source language by running the compiler on 16 programs with interleaved if- and for-statements; these include classical MPC benchmarks such as PSI and Biometric matching, as well as kMeans, Histogram, and other examples from the literature. Our compiler takes the routine and generated *non-vectorized* MOTION code (from MPC Source on the picture above). It then optimizes MPC Source and generates *vectorized* MOTION code (from Optimized MPC Source). We then run the two versions using Boolean GMW and the BMR protocols. (MOTION, which is designed for protocol mixing, supports Arithmetic GMW, however, it does not implement Comparison (CMP) and Multiplexing (MUX) as they would be rather inefficient.) In our experiments vectorized code exhibits 24x improvement on average in circuit generation time for Boolean GMW (20x for BMR), 7x reduction in communication (2x), 97x reduction in number of gates (91x), 4x improvement in setup time in the LAN setting (23x) and 21x improvement in online time, again in the LAN setting (18x).

The compiler can be found at <https://github.com/milana2/ParallelizationForMPC> where <https://github.com/milana2/ParallelizationForMPC/tree/gh-pages> shows experiments on small input sizes and details the different stages of the compiler.

Our results emphasize the importance of backend-independent optimizations – vectorization (described in this work) and protocol mixing (tackled in previous works [14, 25, 21]) are two optimizations readily available at the level of MPC Source. We believe that our work can lead to future work on backend-independent compilation and optimization, ushering new MPC optimizations and combinations of optimizations in the vein of standard compilers,

and thus bringing MPC programming technology closer to practice and wider applicability.

1.2 Outline

The rest of the paper is organized as follows. §2 presents an overview of the compiler. §3 describes our model for cost estimation and argues NP-hardness of optimal scheduling. §4 details the front-end phases of the compiler, §5 focuses in on backend-independent vectorization, and §6 describes translation into MOTION. §7 presents the experimental evaluation. §8 discusses related work and §9 concludes.

All our code, including benchmark Python-like code, compilation phases, and generated MOTION code is available on Github. We do not provide the link only to preserve anonymity, however, we will gladly make it available upon request from reviewers. The Github setup generates graphs and intermediate code for each program along each compiler phase and runs MOTION on small inputs to generate tables of data (the experiments in the paper run on real LAN and WAN). We plan to release the link and code, which we believe will be useful to researchers.

2 OVERVIEW

2.1 Source

As a running example, consider Biometric matching, a standard MPC benchmark. An intuitive (and naive) implementation is as shown in Listing 1(a). Array C is the feature vector of D features that we wish to match and S is the database of N size-D vectors that we match against.

Our compiler takes essentially standard IMP [32] syntax and imposes certain semantic restrictions. (We detail the restrictions in the following sections.) The programmer writes an iterative program and annotates certain inputs and outputs as *shared*. In the example arrays C and S are shared, meaning that they store shares, however, the array sizes D and N respectively are plaintext. The code iterates over the entries in the database and computes the Euclidean distance of the current entry $S[i]$ and C (its square actually). The program returns the index of the vector that gives the best match plus the corresponding sum of squares.

2.2 MPC Source and Cost of Schedule

Our compiler generates an intermediate representation, MPC Source. MPC Source is a *linear* SSA form. MPC Source for Biometric Matching is shown and described in detail in Listing 1(b).

We turn to our analytical model to compute the *cost* of the iterative program. Assume cost β for a local MPC operation (e.g., ADD in Arithmetic sharing) and cost α for a remote MPC operation (e.g., MUX, CMP, etc.). Assuming that ADD is β and SUM, CMP and MUX are α , the MPC Source in Listing 1(b) gives rise to an iterative schedule with cost $ND(2\alpha + \beta) + N(3\alpha)$.

A key contribution is the vectorizing transformation. We can compute all $N * D$ subtraction operations (line 9 in (b)) in a single SIMD instruction; similarly we can compute all multiplication operations (line 10) in a single SIMD instruction. And while computation of an individual sum remains sequential, we can compute the N sums in parallel.

```

233
234
235
236
237
238 1 def biometric(C: shared[list[int]], D: int,
239 2   S: shared[list[int]], N: int) ->
240 3   shared[tuple[int,int]]:
241 4   min_sum : int = MAX_INT
242 5   min_idx : int = 0
243 6   for i in range(N):
244 7     sum : int = 0
245 8     for j in range(D):
246 9       # d = S[i,j] - C[j]
247 10      d : int = S[i * D + j] - C[j]
248 11      p : int = d * d
249 12      sum = sum + p
250 13      if sum < min_sum:
251 14        min_sum : int = sum
252 15        min_idx : int = i
253 16 return (min_sum, min_idx)
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

```

(a) IMP Source

(b) MPC Source

(c) Optimized MPC Source

Table 1: Biometric Matching: ===== From (a) IMP Source to (b) MPC Source: First, MPC Source is an SSA form. Second, it is linear. The conditional in lines 13-14 in IMP Source turns into the linear code in lines 12-16 in MPC Source. The test turns into the CMP operation $t = \text{CMP}(\text{sum!3}, \text{min_sum!2})$, followed by the true-branch sequence, followed by the MUX operations. The first MUX operation selects the value of min_sum : if t is true, then min_sum gets the value of the second multiplexer argument, min_sum!3 , otherwise it takes the value of the third argument, min_sum!2 . Third, MPC Source is a special form of SSA. The SSA ϕ -nodes at the if-then-else (lines 13-15) turn into MUX operations, while the ϕ -nodes at for-loops turn into *pseudo* PHI nodes with a straightforward semantics. ===== From (b) MPC Source to (c) Optimized MPC Source: The compiler determines that SUB and MUL in “naive” MPC Source (lines 9 and 10 in (b)) can be fully vectorized into the SIMD SUB and MUL in optimized MPC Source (lines 9 and 10 in (c)). Notation $p[l,j]$ denotes a 2-dimensional array with fully vectorized dimensions. The computation of sum (line 11 in (b)) is sequential across the j -dimension, but it is parallel across the i -dimension. The loop in lines 12-16 in (c) illustrates; here $p[l,j]$ refers to the j -th column in p . Unfortunately, CMP and MUX remain sequential.

2.3 Vectorized MPC Source and Cost of Schedule

Our compiler produces the vectorized program shown and described in Listing 1(c). Note that this is still our intermediate representation, Optimized MPC Source. Subsequently, the compiler turns this code into MOTION variables, loops and SIMD primitives, which MOTION then uses to generate the circuit.

In MPC back ends, executing n operations “at once” in a single SIMD operation costs a lot less than executing those n operations one by one. This is particularly important when there is communication (i.e., in remote), since many 1-bit values are sent at once rather than sequentially. We elaborate on the cost model in Section §3 but for now consider that each operation has a *fixed* portion (does benefits from amortization) and a *variable* portion (does not benefit from amortization): $\alpha = \alpha_{fix} + \alpha_{var}$. This gives rise to the

following formula for amortized cost: $f(n) = \alpha_{fix} + n\alpha_{var}$, as opposed to unamortized cost $g(n) = n\alpha_{fix} + n\alpha_{var}$. (We extend the same reasoning to β -instructions.)

Thus, the fixed cost of the vectorized program amounts to $2\alpha_{fix} + D\beta_{fix} + N(3\alpha_{fix})$. (The variable cost is the same in both the vectorized and non-vectorized programs.) The first term in the sum corresponds to the vectorized subtraction and multiplication (lines 9-10 in (c)), the second term corresponds to the for-loop on j (lines 12-16) and the third one corresponds to the remaining for-loops on i (lines 19-25). Clearly, $2\alpha_{fix} + D\beta_{fix} + N(3\alpha_{fix}) \ll ND(2\alpha_{fix} + \beta_{fix}) + N3\alpha_{fix}$. Empirically, we observe that $\alpha_{var} \approx 0$ and orders of magnitude improvement (e.g., 12x in online time in GMW). Additionally, the unvectorized version runs out of memory for $N = 256$, while the vectorized one runs with the standard maximal input size $N = 4,096$.

3 ANALYTICAL MODEL

This section presents a model to reason about the cost of execution of MPC programs, including accounting for amortization. We define the assumptions and setting in §3.1. We proceed to define the scheduling problem in §3.2, which we expected to be able to solve optimally. §3.3 shows that the problem is NP-hard via a reduction to the Shortest Common Supersequence (SCS) problem. Despite the negative general result, we expect the formulation in terms of SCS to be useful as sequences are short and few in practice.

3.1 Scheduling in MPC

For this treatment we make the following simplifying assumptions:

- (1) All statements in the program execute using the same protocol (sharing). That is, there is no share conversion.
- (2) There are two tiers of MPC instructions, local and remote. A local instruction (e.g., ADD in Arithmetic, XOR in Boolean) has cost β and a remote instruction (e.g., MUX, MUL, SHL, etc.) has cost α , where $\alpha \gg \beta$. We assume that all remote instructions have the same cost.
- (3) In MPC frameworks, executing n operations “at once” in a single SIMD operation costs a lot less than executing those n operations one by one. Following Amdahl’s law, we write $\alpha = \frac{1}{s}p\alpha + (1-p)\alpha$, where p is the fraction of execution time that benefits from amortization and $(1-p)$ is the fraction that does not, and s is the available resource. Thus, $n\alpha = \frac{n}{s}p\alpha + n(1-p)\alpha$. For the purpose of the model we assume that s is large enough and the term $\frac{n}{s}p\alpha$ amounts to a *fixed cost* incurred regardless of whether n is 10,000 or just 1. (This models the cost of preparing and sending a packet from party A to party B for example.) Therefore, amortized execution of n operations is $f(n) = \alpha_{fix} + n\alpha_{var}$ in contrast to unamortized execution $g(n) = n\alpha_{fix} + n\alpha_{var}$. We have $\alpha_{fix} \ll n\alpha_{fix}$ and since fixed cost dominates variable cost (particularly for remote operations), we have $f(n) \ll g(n)$.
- (4) MPC instructions scheduled in parallel benefit from amortization *only if* they are the same instruction. Given our previous assumption, 2 MUL instructions can be amortized in a single SIMD instruction that costs $\alpha_{fix} + 2\alpha_{var}$, however a MUL and a MUX instruction still cost $2\alpha_{fix} + 2\alpha_{var}$ even when scheduled “in parallel”.¹

3.2 Problem Statement

As mentioned earlier, at the lowest level, we have two types of MPC instructions (also called *gates* or *operations* in similar works) 1) local/non-interactive (e.g., an addition instruction A) and 2) remote/interactive (e.g., a multiplication instruction M).

Given a serial schedule (a linear graph) of an MPC program i.e. a sequence of instructions $S := (S_1; \dots; S_n)$, where $S_i \in \{A, M\}$, $1 \leq i \leq n$, and a def-use dependency graph $G(V, E)$ corresponding to S , our task is to construct a parallel schedule (another linear graph) $P := (P_1; \dots; P_m)$ observing the following conditions:

- (1) All P_i ’s consist of MPC instructions of the same kind, e.g., all MUL, MUX, ADD, etc.

¹This is not strictly true, but assuming it, e.g. as in [25, 18, 30], helps simplify the problem.

- (2) Def-use dependencies of the graph $G(V, E)$ are preserved i.e. if instructions $S_i, S_j, i < j$ form a def-use i.e. an edge exists from S_i to S_j in G , then they can only be mapped to $P_{i'}, P_{j'}$ such that $i' < j'$.

Correctness. Correctness of P is guaranteed by definition. Preserving def-use *dependencies* means the computed function remains the same in both S and P .

The cost of schedule S is

$$\text{cost}(S) = \sum_{i=1}^n \text{cost}(S_i) = L_\alpha \alpha_{fix} + L_\beta \beta_{fix} + L_\alpha \alpha_{var} + L_\beta \beta_{var} \quad (1)$$

where L_α is the number of α -instructions and L_β is the number of β ones. (We used this formula to compute the cost of the unrolled MPC Source program in §2.) The cost of schedule P is more interesting:

$$\text{cost}(P) = \sum_{i=1}^m \text{cost}(P_i) \quad (2)$$

Each P_i may contain multiple instructions, and $\text{cost}(P_i)$ is amortized. Thus, according to our model $\text{cost}(P_i) = \alpha_{fix} + |P_i| \alpha_{var}$ if P_i stores $|P_i|$ α -instructions, or $\text{cost}(P_i) = \beta_{fix} + |P_i| \beta_{var}$ if it stores β -instructions. (Similarly, we used this formula to compute the cost of the Optimized MPC Source program in §2.)

Our goal is to construct a parallel schedule P that reduces the program cost (when compared to cost of S), possibly an optimal schedule. Originally we hoped that the problem is simpler and computation of the optimal schedule is tractable. Unfortunately, the optimal schedule turns out to be NP-hard via a reduction to the Shortest Common Supersequence problem.

3.3 Scheduling is NP-hard

To prove that optimal scheduling is an NP-Hard problem, we consider the following convenient representation. An MPC program is represented as a set of sequences $\{s_1, \dots, s_n\}$ of operations. In each sequence s_j operations depend on previous operations via a def-use i.e. $s_i[j], j > 1$ depends on $s_i[j-1]$.

As an example, consider the MPC program consisting of the following three sequences, all made up of two distinct α -instructions M_1 and M_2 , e.g., M_1 is MUL and M_2 is MUX. The right arrow indicates a def-use *dependence*, meaning that the source node must execute before the target node:

- (1) $M_1 \rightarrow M_2 \rightarrow M_1$
- (2) $M_1 \rightarrow M_1 \rightarrow M_1$
- (3) $M_2 \rightarrow M_1 \rightarrow M_2$

The problem is to find a schedule P with *minimal cost*. For example, a schedule with minimal cost for the sequences above is

$$M_1(1), M_1(2); M_1(2); M_2(1), M_2(3); M_1(1), M_1(2), M_1(3); M_2(3)$$

The parentheses above indicate the sequence where the instruction comes from: (1), (2), or (3). Cost of schedule P is computed using Eq. (2) above and it amounts to $5\alpha_{fix} + 9\alpha_{var}$.

The problem of finding a schedule P with a minimal $\text{cost}(P)$ is shown to be NP-Hard problem, as it can be reduced to the problem of finding a *shortest common supersequence*, a known NP-Hard problem [38]. The shortest common supersequence problem is as follows: *given two or more sequences find the the shortest sequence*

that contains all of the original sequences. This can be solved in $O(n^k)$ time, where n is the cardinality of the longest sequence and k is the number of sequences. We can see that the optimal schedule is the shortest schedule, since the shortest schedule minimizes the fixed cost while the variable cost remains the same.

To formalize the reduction, suppose P is a schedule with minimal cost (computed by a black-box algorithm). Clearly P is a supersequence of each sequence s_i , i.e., P is a common supersequence of $s_1 \dots s_n$. It is also a shortest common supersequence. The cost of $\text{cost}(P) = L\alpha_{\text{fix}} + N\alpha_{\text{var}}$ where L is the length of P and N is the total number of instructions across all sequences. Now suppose, there exist a shorter common supersequence P' of length L' . $\text{cost}(P') < \text{cost}(P)$ since $L'\alpha_{\text{var}} + N\alpha_{\text{var}} < L\alpha_{\text{var}} + N\alpha_{\text{var}}$, contradicting the assumption that P has the lowest cost. \square

4 COMPILER FRONT END

Fig. 1 presents an overview of our compiler. We start the section with a description of the source syntax and semantic restrictions in §4.1. We proceed to describe the front end of the compiler: §4.2 describes translation from IMP Source to SSA, and §4.3 describes translation from SSA into MPC Source. We detail analysis on MPC Source and backend-independent vectorization in §5 and translation into MOTION code in §6.

A key novelty of our work is the systematic translation of high-level source, which includes if-then-else statements, into a linear sequence of primitive MPC instructions. We start with an AST source syntax, then convert it into a Control-flow graph (CFG), as Cytron's classical SSA algorithm [16] is defined over a CFG. SSA is the natural means for converting if-then-else statements into MUX primitives, as there is correspondence between phi-nodes and MUX primitives (§4.2). Yet, to the best of our knowledge, we are the first to present an algorithm that uses Cytron's SSA conversion. Following translation into an SSA CFG, we then translate into MPC Source, which is most conveniently described with an AST syntax (§4.3).

We also note that we settled on this process, going from AST to CFG then back to an AST representation, after the straightforward approach failed. Originally, we attempted to reuse existing implementations of SSA intermediate representations such as Soot Shimple (<http://soot-oss.github.io/soot/>) or LLVM IR (<https://llvm.org/>). The problem was that these were CFG representations and they had lost connection between the ϕ node and the conditional that triggered the ϕ node. Specifically, given $x = \phi(y, z)$, there is no information what conditional triggered the ϕ node and whether y corresponds to the true or false branch of the conditional. Moreover, ϕ nodes with 3 or more arguments are common. Yet another disadvantage of existing SSA IRs was that they did not handle arrays, not even the handling proposed in Cytron's work [16]. Handling of arrays is important for optimizations such as vectorization. And while it is possible to reconstruct the missing information via some form of control dependence analysis, and it is possible to add arrays, this proved very difficult due to the complexity of the IRs, as they are designed to handle much richer and complex syntax and semantics than MPC. A clean state solution, where we started from the AST and retained all necessary information in the CFG,

proved the correct choice and drove our progress on the rest of the compiler.

4.1 Syntax and Semantic Restrictions

Source syntax is standard IMP syntax but with for-loops:

```

e ::= e op e | x | const | A[e]  expression
s ::= s; s |                      sequence
x = e | A[e] = e |                assignment stmt
for i in range(I) : s |           for stmt
if e: s else: s                    if stmt
    
```

The syntax allows for array accesses, arbitrarily nested loops, and if-then-else control flow.

The prototype version of the compiler assumes the following semantic restrictions on source programs. Currently, the compiler does not enforce the restrictions, however, they can be easily encoded as attribute rules, also known as syntax-directed translation [2, 35], over the AST syntax above. The reason why we do not implement the rules is because the majority of these restrictions are implementation restrictions that can be lifted in future versions of our compiler.

- (1) Loops are of the form $0 \leq i < I$ and bounds are fixed at compile time. It is a standard restriction in MPC that the bounds must be known at circuit-generation time.
- (2) Arrays are one-dimensional. N -dimensional arrays are linearized and accessed in row-major order and at this point the programmer is responsible for linearization and access. This restriction can be easily lifted.
- (3) Array subscripts are plaintext values. This restriction can be easily lifted by applying the standard linear scan [29, 5] when the subscript is a secret-shared value.
- (4) The subscript e is a function of the indices of the enclosing loops. For read access, the compiler allows an arbitrary such function. However, it restricts write access to *canonical* writes, i.e., $A[i, j, k] = \dots$ where i, j and k are the indices of the *outermost* loops enclosing the array write statement. These indices loop over the three dimensions of A and all write access to A follow this. This restriction simplifies vectorization in the presence of array writes. We plan to extend the compiler with arbitrary write access.
- (5) The final restriction disallows array writes from within if-then-else statements. This is to ensure that arguments of MUX in the MPC Source translation are base types, i.e., just int or bool, and not array types. Again, this restriction simplifies our current implementation and it will be lifted, though not as trivially as the first two. The restriction causes a minor inconvenience to the programmer as they may not write

```
1 if e: A[i] = val
```

Instead they write

```
1 if not(e): val = A[i]
2 A[i] = val
```

In addition, our compiler defines and implements a taint type system at the level of MPC Source. We define the base MPC Source syntax and the type system in §4.4. We note that while the programmer writes annotations at the level of IMP Source (as in Listing 1(a)),

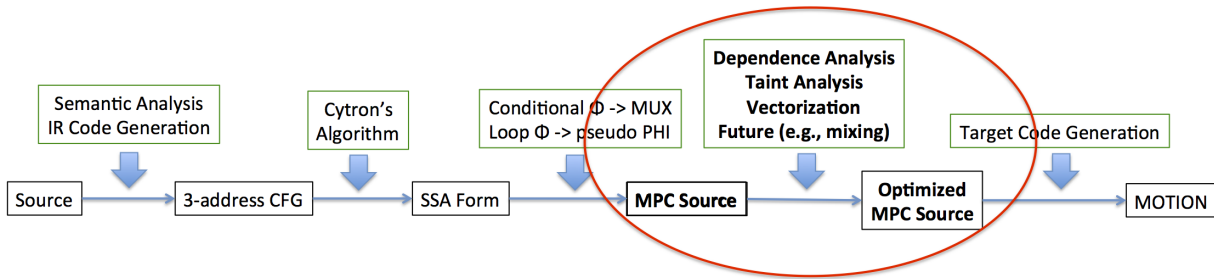


Figure 1: Compiler Framework.

the annotations propagate through the transformations; annotations are inferred and checked with a taint analysis (based on the type system) at the level of MPC Source. The only required annotations are on the input arguments.

For the rest of this section we write i, j, k to denote the loop nest: i is the outermost loop, j , is immediately nested in i , and so on until k and we use I, J, K to denote the corresponding upper bounds. We write $A[i, j, k]$ to denote canonical access to an array element. In the program, canonical access is achieved via the standard row-major order formula: $(J * K) * i + K * j + k$. To simplify the presentation we describe our algorithms in terms of three-element tuples i, j, k , however, discussion easily generalizes to arbitrarily large loop nests.

4.2 From IMP Source to SSA

Our compiler translates from Source to SSA as follows:

Parsing: Use Python’s `ast` module to parse the input source code to a Python AST.

Syntax checking: Ensure that the AST matches the restricted subset defined in Section §4.1. This step outputs an instance of the `restricted_ast`. Function class, which represents our restricted subset of the Python AST.

3-address CFG conversion: Convert the restricted-syntax AST to a three-address control-flow graph. To do this, first, add an empty basic block to the CFG and mark it as current. Next, for each statement in the restricted AST’s function body, process the statement. Statements can either be for-loops, if-statements, or assignments (as in §4.1). Rules for processing each kind of statement are given below:

- (1) **For-loops:** Create new basic blocks for the loop condition (the *condition-block*), the loop body (the *body-block*), and the code after the loop (the *after-block*). Insert a jump from the end of the current block to the condition-block. Then, mark the condition-block as the current block. Insert a for-instruction at the end of the current block with the loop counter variable and bounds from the AST. Next, add an edge from the current block to the after-block labeled “FALSE” and an edge from the current block to the body-block labeled “TRUE”. Then, set the body-block to be the current block and process all statements in the AST’s loop body. Finally, insert a jump to the condition-block and set the after-block as current.

- (2) **If-statements:** Create new basic blocks for the “then” statements of the if-statement (the *then-block*), the “else” statements of the if-statement (the *else-block*), and the code after the if-statement (the *after-block*). At the end of the current block, insert a conditional jump to the then-block or else-block depending on the if-statement condition in the AST. Next, mark the then-block as current, process all then-statements, and add a jump to the after-block. Similarly, mark the else-block as current, process all else-statements, and add a jump to the after-block. Finally, set the after-block to be the current block, and give it a *merge condition* property equal to the condition of the if-statement.

- (3) **Assignments:** In the restricted-syntax AST, the left-hand side of assignments can be a variable or an array subscript. If it is an array subscript, e.g., $A[i] = x$, change the statement to $A = \text{Update}(A, i, x)$. If the statement is not already three-address code, for each sub-expression in the right-hand side of the assignment, insert an assignment to a temporary variable.

SSA conversion: Convert the 3-address CFG to SSA with Cytron’s algorithm.

4.3 From SSA to MPC Source

Once the compiler converts the code to SSA, it transforms ϕ -nodes that correspond to if-statements into MUX nodes. From the 3-address CFG conversion step, ϕ -nodes corresponding to if-statements will be in a basic block with the merge condition property. For example, if $X!3 = \phi(X!1, X!2)$ is in a block with merge condition C , the compiler transforms it into $X!3 = \text{MUX}(C, X!1, X!2)$. Next, the compiler runs the dead code elimination algorithm from Cytron’s SSA paper.

Next, the control-flow graph is *linearized* into MPC Source, which has loops but no if-then-else-statements. This means that both branches of all if-statements are executed, and the MUX nodes determine whether to use results from the then-block or from the else-block. The compiler linearizes the control-flow graph with a variation of depth-first search. Blocks with the “merge condition” property are only considered the second time they are visited, since that will be after both branches of the if-statement are visited. (The Python AST naturally gives rise to a translation where each conditional has exactly two targets, and each “merge condition” block has exactly two incoming edges, a TRUE and a FALSE edge. Thus, each ϕ -node has exactly two multiplexer arguments, which dovetails

into MUX. This is in contrast with Cytron's algorithm which operates at the level of the CFG and allows for ϕ -nodes with multiple arguments.) Each time the compiler visits a block, it adds the block's statements to the MPC source. If the block ends in a for-instruction, the compiler recursively converts the body and code after the loop to MPC source and adds the for-loop and code after the loop to the main MPC source. If the block does not end in a for-instruction, the compiler recursively converts all successor branches to MPC source and appends these to the main MPC source.

```

{ Step 1: Replace  $\phi$ -nodes with MUX nodes }
for each basic block block in the control-flow graph do
  if block has the merge condition property then
    merge_condition  $\leftarrow$  merge condition variable of block
    for each  $\phi$ -node phi =  $\phi(v_1, v_2)$  in block do
      Replace phi with  $MUX(\text{merge\_condition}, v_1, v_2)$  in block
    end for
  end if
end for
{ Step 2: Linearize control-flow graph into MPC Source }
visited  $\leftarrow$  empty set
merge_visited  $\leftarrow$  empty set
Define search(block):
  if block has the merge condition property and block is not in
  merge_visited then
    Add block to merge_visited
  return empty list
end if
Add block to visited
if block is a for-loop header then
  cfg_body  $\leftarrow$  successor of block containing the beginning of
  the for-loop body
  cfg_after  $\leftarrow$  successor of block containing the beginning of
  the code after the for-loop
  mpc_body  $\leftarrow$  list of the  $\phi$ -functions in block concatenated
  with search(cfg_body)
  loop  $\leftarrow$  MPC Source for-loop statement with the same
  counter variable and bounds as block and with mpc_body
  as its body
  return loop prepended to search(cfg_after)
else
  result  $\leftarrow$  empty list
  for each successor successor of block do
    if successor is not in visited then
      result  $\leftarrow$  result concatenated with search(successor)
    end if
  end for
  return result
end if
return search(entry block of the control-flow graph)

```

Now, the remaining ϕ -nodes in MPC source are the loop header nodes. We call these nodes *pseudo* ϕ -nodes and we write PHI in MPC Source. A pseudo ϕ -node $x!1 = \text{PHI}(x!0, x!2)$ in a loop header is evaluated during circuit generation. If it is the 0-th iteration, then the ϕ -node evaluates to $x!0$, otherwise, it evaluates to $x!2$.

4.4 Base MPC Source Syntax and Taint Types

The syntax of the MPC Source program produced by the above section is essentially IMP syntax. (In section 5 we extend the base syntax to account for vectorization.) Most notably, there is no if-then-else statement but there are MUX expressions:

$$\begin{array}{ll}
 e ::= e \text{ op } e \mid x \mid \text{const} \mid A[e] \mid \text{MUX}(e, e, e) & \text{expression} \\
 s ::= s; s \mid & \text{sequence} \\
 x = \text{PHI}(x, x) \mid x = e \mid A[e] = e \mid & \text{assignment stmt} \\
 \text{for } i \text{ in range}(I) : s & \text{for stmt}
 \end{array}$$

Expressions are typed $\langle q \tau \rangle$, where qualifier q and type τ are:

$$\begin{array}{ll}
 \tau ::= \text{int} \mid \text{bool} \mid \text{list}[\text{int}] \mid \text{list}[\text{bool}] & \text{base types} \\
 q ::= \text{shared} \mid \text{plain} & \text{qualifiers}
 \end{array}$$

The type system is standard, and in our experience, a sweet spot between readability and expressivity. The shared qualifier denotes shared values, i.e., ones shared among the parties and computed upon under secure computation protocols; the plain qualifier denotes plaintext values. Subtyping is $\text{plain} <: \text{shared}$, meaning that we can convert a plaintext value into a shared one, but not vice versa. Subtyping on qualified types is again as expected, it is covariant in the qualifier and invariant in the type: $\langle q_1 \tau_1 \rangle <: \langle q_2 \tau_2 \rangle$ iff $q_1 <: q_2$ and $\tau_1 = \tau_2$.

The typing rules for non-trivial expressions are as follows:

$$\frac{\Gamma \vdash e_1 : \langle q_1 \tau \rangle \quad \Gamma \vdash e_2 : \langle q_2 \tau \rangle \quad \tau \in \{\text{int}, \text{bool}\}}{\Gamma \vdash e_1 \text{ op } e_2 : \langle q_1 \vee q_2 \tau \rangle} \text{(OP)}$$

$$\frac{\Gamma \vdash e : \langle \text{plain int} \rangle \quad \Gamma \vdash A : \langle q \text{ list}[\tau] \rangle \quad \tau \in \{\text{int}, \text{bool}\}}{\Gamma \vdash A[e] : \langle q \tau \rangle} \text{(ARRAY ACCESS)}$$

$$\frac{\Gamma \vdash e_1 : \langle q_1 \text{ bool} \rangle \quad \Gamma \vdash e_2 : \langle q_2 \tau \rangle \quad \Gamma \vdash e_3 : \langle q_3 \tau \rangle \quad \tau \in \{\text{int}, \text{bool}\}}{\Gamma \vdash \text{MUX}(e_1, e_2, e_3) : \langle q_1 \vee q_2 \vee q_3 \tau \rangle} \text{(MUX)}$$

Similarly, the typing rules for statements are as follows. The constraints are standard: the right-hand side of an assignment is a subtype of the left-hand side.

$$\frac{\Gamma \vdash x_1 : \langle q_1 \tau \rangle \quad \Gamma \vdash x_2 : \langle q_2 \tau \rangle \quad \Gamma \vdash x_3 : \langle q_1 \tau \rangle \quad q_2 \vee q_3 <: q_1}{\Gamma \vdash x_1 = \text{PHI}(x_2, x_3) : \text{OK}} \text{(PHI ASSIGN)}$$

$$\frac{\Gamma \vdash x : \langle q_1 \tau \rangle \quad \Gamma \vdash e : \langle q_2 \tau \rangle \quad q_2 <: q_1 \quad \tau \in \{\text{int}, \text{bool}\}}{\Gamma \vdash x = e : \text{OK}} \text{(VAR ASSIGN)}$$

$$\frac{\Gamma \vdash i : \langle \text{plain int} \rangle \quad \Gamma \vdash I : \langle \text{plain int} \rangle \quad \Gamma \vdash s : \text{OK}}{\Gamma \vdash \text{for } i \text{ in range}(I) : s : \text{OK}} \text{(FOR STMT)}$$

As mentioned earlier, the only annotations the program need provide is on program inputs. The compiler infers the rest of the annotations. The type system has two purposes (1) it imposes restrictions, and (2) it enables code generation, specifically, it informs the backend on whether a statement operates on shared variables or plaintext ones, and the backend generates appropriate MOTION code.

5 BACKEND-INDEPENDENT VECTORIZATION

This section describes our vectorization algorithm. While vectorization is a longstanding problem, and we build upon existing work on scalar expansion and classical loop vectorization [4], our algorithm is unique as it works on the MPC Source SSA-form representation. We posit that vectorization over MPC Source is a new problem that warrants a new look, in part because of MPC's unique linear structure and in part because vectorization meshes in with other MPC-specific optimizations in non-trivial ways (other works have explored manual vectorization and protocol mixing in an ad-hoc way, e.g., [18, 13, 25]).

§5.1 describes our dependence analysis and §5.2 describes scalar expansion, which lifts scalars (and arrays) to the corresponding loop dimensionality to create opportunities for vectorization. §5.3 describes our core vectorization algorithm and §5.4 argues correctness of the vectorization transformation. §5.5 extends vectorization with array writes.

5.1 Dependence Analysis

We build a dependence graph where the nodes are the MPC Source statements and the edges represent the def-use relations.

Def-use Edges. We distinguish the following def-use edges:

- same-level edge $X \rightarrow Y$ where X and Y are in the same loop nest, say i, j, k . E.g., the def-use edge 9 to 10 in the Biometric MPC Source in Listing 1 is a same-level edge. A same-level edge can be a back-edge in which case a ϕ node is the target of the edge. E.g., 15 to 4 in Biometric is a same-level back-edge.
- outer-to-inner $X \rightarrow Y$ where X is in an outer loop nest, say i , and Y is in an inner one, say i, j, k . E.g., 1 to 4 in Biometric forms is an outer-to-inner edge.
- inner-to-outer $X \rightarrow Y$ where X is a *phi*-node in an inner loop nest, i, j, k , and Y is in the enclosing loop nest i, j . E.g., the def-use from 8 to 12 gives rise to an inner-to-outer edge. An inner-to-outer edge can be a back-edge as well, in which case both X and Y are ϕ -nodes with the source X in a loop nested into Y 's loop (not necessarily immediately).
- mixed forward edge $X \rightarrow Y$. X is in some loop i, j, k and Y is in a loop nested into i, j, k' . We transform mixed forward edges as follows. Let x be the variable defined at X . We add a variable and assignment $x' = x$ immediately after the i, j, k loop. Then we replace the use of x at Y with x' . This transforms a mixed forward edge into an "inner-to-outer" forward edge followed by an outer-to-inner forward edge. Thus, Basic Vectorization handles one of "same-level", "inner-to-outer", or "outer-to-inner" def-use edges.

Closures. We define $\text{closure}(n)$ where n is a PHI-node. Intuitively, it computes the set of nodes (i.e., statements) that form a dependence cycle with n . The closure of n is defined as follows:

- n is in $\text{closure}(n)$
- X is in $\text{closure}(n)$ if there is a same-level path from n to X , and $X \rightarrow n$ is a same-level back-edge.
- Y is in $\text{closure}(n)$ if there is a same-level path from n to Y and there is a same-level path from Y to some X in $\text{closure}(n)$.

5.2 Scalar Expansion

An important component of our algorithm is the scalar expansion to the corresponding loop dimensionality, which is necessary to expose opportunities for vectorization. In the Biometric example, $d = S[i * D + j] - C[j]$ equiv. to $d = S[i, j] - C[j]$, which gave rise to $N * D$ subtraction operations in the sequential schedule, is lifted. The argument arrays S and C are lifted and the scalar d is lifted: $d[i, j] = S[i, j] - C[j]$. The algorithm then detects that the statement can be vectorized.

Arrays. Conceptually, we treat all variables as arrays. There are three kinds of arrays.

- Scalars: These are scalar variables we expand into arrays for the purposes of vectorization. For those, all writes are canonical writes and all reads are canonical reads. We will *raise dimension* when a scalar gives rise to an outer-to-inner dependence edge (e.g., `sum2` in line 6 of the MPC source code will be raised to a 1-dimensional array since `sum2` is used in the inner j -loop). We will *drop dimension* when a scalar gives rise to an inner-to-outer dependence edge (e.g., `sum2` for which the lifted inner loop computes D values, but the outer loop only needs the last one.)
- Read-only input arrays: There are no writes, while we may have non-canonical reads, $f(i, j, k)$. Vectorization adds raise dimension operations at the beginning of the function to lift these arrays to the dimensionality of the loop where they are used, possibly *reshaping* the arrays. If there are multiple "views" of the input array, there would be multiple raise dimension statements to create each one of these views. The invariant is that at reads in loops, the reads of "views" of the original input array are canonical.
- Read-write output arrays: Writes are canonical (by restriction) but reads can be non-canonical. Dependences limit vectorization when non-canonical read access refers to array writes in previous iterations, thus creating loop-carried dependences. We may apply both raise and drop dimension, however, they respect the fixed dimensionality of the output array. The array cannot be raised to a dimension lower than its canonical (fixed) dimensionality and it cannot be dropped to lower dimension. In addition, non-canonical reads may require lifting (i.e., reshaping) of the array after the most recent write rather than in the beginning of the program in order to reduce a non-canonical read to a canonical one.

Raise dimension. The `raise_dim` function expands a scalar (or array). There are two conceptual versions of `raise_dim`. One reshapes an arbitrary access into a canonical read access in the corresponding loop. It takes the original array, the access pattern function $f(i, j, k)$

in loop nest i, j, k and the loop bounds $((i : I), (j : J), (k : K))$:

$$\text{raise_dim}(A, f(i, j, k), ((i : I), (j : J), (k : K)))$$

It produces a new 3-dimensional array A' by iterating over i, j, k and setting each element of A' as follows:

$$A'[i, j, k] = A[f(i, j, k)]$$

The end result is that uses of $A[f(i, j, k)]$ in loop nest i, j, k are replaced with canonical read-accesses to $A'[i, j, k]$ that can be vectorized. In the running Biometric example, $C' = \text{raise_dim}(C, j, (i : N, j : D))$ lifts the 1-dimensional array C into a 2-dimensional array. The i, j loop now accesses C' in the canonical way, $C'[i, j]$. Similarly, $S' = \text{raise_dim}(S, i * D + j, (i : N, j : D))$ tries to lift S , but the operation turns into a no-op because S is already a 2-dimensional array and the read access is canonical.

The other version of raise_dim lifts a lower-dimension array into a higher-dimension for access in a nested loop. It is necessary when processing outer-to-inner dependences. Here A is an i -array and raise dimension adds two additional dimensions:

$$\text{raise_dim}(A, (j : J, k : K))$$

This version is reduced to the above version by adding the access pattern function, which is just i :

$$\text{raise_dim}(A, i, (j : J, k : K))$$

Drop dimension. The corresponding drop_dim is carried out when an array written in an inner loop is used in an enclosing loop. It takes a higher dimensional array, say i, j, k and removes trailing dimensions, say j, k :

$$\text{drop_dim}(A, (j : J, k : K))$$

It iterates over i and takes the result at the maximal index of j and k , i.e., the result at the last iterations of j and k :

$$A'[i] = A[i, J - 1, K - 1]$$

5.3 Basic Vectorization

We present our algorithm followed by an example on Biometric.

5.3.1 Algorithm

There are two key phases of the algorithm. Phase 1 inserts raise dimension and drop dimension operations according to def-uses. E.g., if there is an inner-to-outer dependence, it inserts raise_dim , and similarly, if there is an outer-to-inner dependence, it inserts drop_dim . After this phase operations work on arrays of the corresponding dimensionality and we optimistically vectorize all arrays.

Phase 2 proceeds from the inner-most towards the outer-most loop. For each loop it anchors dependence cycles (closures) around pseudo PHI nodes then removes vectorization from the dimension of that loop. There are two important points in this phase. First, it may break a loop into smaller loops which would discover opportunities for vectorization in intermediate statements in the loop. Second, it handles writes arrays. It creates opportunities for vectorization in the presence of write arrays, even though Cytron's SSA adds a back-edge to the array PHI-node, thus killing vectorization of statements that read and write that array.

The excerpts in red color in the pseudo code below highlight the extension with array writes. We advise the reader to omit the extension for now and consider just read-only arrays. We explain

the extension in §5.5. (As many of our benchmarks include write arrays, it plays an important role.)

Phase 3 cleans up local arrays of references (this is an optional phase and our current implementation does not include it) and Phase 4 explicitly turns operations into MOTION SIMD operations.

{ Phase 1: Raise dimension of scalar variables to corresponding loop nest. We can traverse stmts linearly in MPC-source. }

for each MPC *stmt* : $x = \text{Op}(y_1, y_2)$ in loop i, j, k **do**

for each argument y_n **do**

case def-use edge $\text{stmt}'(\text{def of } y_n) \rightarrow \text{stmt}(\text{def of } x)$ of
same-level: y'_n is y_n

outer-to-inner: add $y'_n[i, j, k] = \text{raise_dim}(y_n)$ at stmt'
(more precisely, right after stmt')

inner-to-outer: add $y'_n[i, j, k] = \text{drop_dim}(y_n)$ at stmt
(more precisely, in loop of stmt right after loop of stmt')

end for

{ Optimistically vectorize all. I means vectorized dimension. }
change to $x[I, J, K] = \text{Op}(y'_1[I, J, K], y'_2[I, J, K])$

end for

{ Phase 2: Recreating for-loops for cycles; vectorizable stmts hoisted up. }

for each dimension d from highest to 0 **do**

for each PHI-node n in loop i_1, \dots, i_d **do**

compute $\text{closure}(n)$

end for

{ cl_1 and cl_2 intersect if they have common statement or update same array; "intersect" definition can be expanded }

while there are closure cl_1 and cl_2 that intersect **do**
merge cl_1 and cl_2

end while

for each closure cl (after merge) **do**

create **for** i_d **in** ... loop

add PHI-nodes in cl to header block

add target-less PHI-node for A if cl updates array A

add statements in cl to loop in some order of dependences
{ Dimension is not vectorizable: }

change I_d to i_d in all statements in loop

treat **for**-loop as monolith node for def-uses: e.g., some def-use edges become same-level.

end for

for each target-less PHI-node $A!1 = \text{PHI}(A!0, A!k)$ **do**

in vectorizable stmts, replace use of $A!1$ with $A!0$

discard PHI-node if not used in any cl , replacing $A!1$ with $A!0$ or $A!k$ appropriately

end for

end for

{ Phase 3: Remove unnecessary dimensionality. }

{ A dimension i is dead on exit from stmt $x[\dots i \dots] = \dots$ if all def-uses with targets outside of the enclosing **for** i ... MOTION loop end at target (use) $x' = \text{drop_dim}(x, i)$. }

for each stmt and dimension $x[\dots i \dots] = \dots$ **do**

if i is a dead dimension on exit from stmt $x[\dots i \dots] = \dots$, remove i from x (all defs and uses)

end for

{ Now clean up drop_dim and raise_dim }

for each $x' = \text{drop_dim}(x, i)$ **do**

```

1045     replace with  $x' = x$  if  $i$  is dead in  $x$ .
1046   end for
1047   do (1) (extended) constant propagation, (2) copy propagation and
1048   (3) dead code elimination to get rid of redundant variables and
1049   raise and drop dimension statements
1050   { Phase 4: }
1051   add SIMD for simdfied dimensions

```

5.3.2 Example: Biometric

We now demonstrate the workings of the basic algorithm on our running example. Recall the MPC Source for Biometric:

```

1056 1 min_sum!1 = MAX_INT
1057 2 min_idx!1 = 0
1058 3 for i in range(0, N):
1059 4   min_sum!2 = PHI(min_sum!1, min_sum!4)
1060 5   min_idx!2 = PHI(min_idx!1, min_idx!4)
1061 6   sum!2 = 0
1062 7   for j in range(0, D):
1063 8     sum!3 = PHI(sum!2, sum!4)
1064 9     d = SUB(S(((i * D) + j)), C[j])
1065 10    p = MUL(d,d)
1066 11    sum!4 = ADD(sum!3,p)
1067 12    t = CMP(sum!3,min_sum!2)
1068 13    min_sum!3 = sum!3
1069 14    min_idx!3 = i
1070 15    min_sum!4 = MUX(t, min_sum!3, min_sum!2)
1071 16    min_idx!4 = MUX(t, min_idx!3, min_idx!2)
1072 17 return (min_sum!2, min_idx!2)

```

Phase 1 of Vectorization Algorithm. The transformation preserves the dependence edges. It raises the dimensions of scalars and optimistically vectorizes all operations. The next phase discovers loop-carried dependences and removes affected vectorization.

In the code below statements (e.g., $\text{min_sum!3} = \text{sum!3}$) are implicitly vectorized. The example illustrates the two different versions of *raise_dim*. E.g., $\text{raise_dim}(C, j, (i:N,j:D))$ reshapes the read-only input array. $\text{drop_dim}(\text{sum!3})$ removes the j dimension of sum!3 .

```

1083 1 min_sum!1 = MAX_INT
1084 2 min_sum!1^ = raise_dim(min_sum!1, (i:N))
1085 3 min_idx!1 = 0
1086 4 min_idx!1^ = raise_dim(min_idx!1, (i:N))
1087 5 S^ = raise_dim(S, ((i * D) + j), (i:N,j:D))
1088 6 C^ = raise_dim(C, j, (i:N,j:D))
1089 7 for i in range(0, N):
1090 8   min_sum!2 = PHI(min_sum!1^, min_sum!4)
1091 9   min_idx!2 = PHI(min_idx!1^, min_idx!4)
1092 10  sum!2 = 0 // Will lift, when hoisted
1093 11  sum!2^ = raise_dim(sum!2, (j:D))
1094 12  for j in range(0, D):
1095 13    sum!3 = PHI(sum!2^, sum!4)
1096 14    d = SUB(S^, C^ )
1097 15    p = MUL(d,d)
1098 16    sum!4 = ADD(sum!3,p)
1099 17    sum!3^ = drop_dim(sum!3)
1100 18    t = CMP(sum!3^,min_sum!2)
1101 19    min_sum!3 = sum!3^

```

```

20   min_idx!3 = i // Same-level, will lift when hoisted
21   min_sum!4 = MUX(t, min_sum!3, min_sum!2)
22   min_idx!4 = MUX(t, min_idx!3, min_idx!2)
23 min_sum!2^ = drop_dim(min_sum!2)
24 min_idx!2^ = drop_dim(min_idx!2)
25 return (min_sum!2^, min_idx!2^)

```

Phase 2 of Vectorization Algorithm. This phase analyzes statements from the innermost loop to the outermost. The key point is to discover loop-carried dependencies and re-introduce loops whenever dependencies make this necessary.

Starting at the inner phi-node $\text{sum!3} = \text{PHI}(\dots)$, the algorithm first computes its closure. The closure amounts to the phi-node itself and the addition node $\text{sum!4} = \text{ADD}(\text{sum!3}, p!3)$, accounting for the loop-carried dependency of the computation of sum . The algorithm replaces this closure with a for-loop on j removing vectorization on j . Note that the SUB and MUL computations remain outside of the loop as they do not depend on PHI-nodes that are part of cycles. The dependences are from $p[!j] = \text{MUL}(d[!j], d[!j])$ to the monolithic for-loop and from the for-loop to $\text{sum!3} \hat{=} \text{drop_dim}(\text{sum!3})$. Lower case index, e.g., i , indicates non-vectorized dimension, while uppercase index, e.g., I indicates vectorized dimension.

After processing the inner loop code becomes:

```

1 min_sum!1 = MAX_INT
2 min_sum!1^ = raise_dim(min_sum!1, (i:N!0))
3 min_idx!1 = 0
4 min_idx!1^ = raise_dim(min_idx!1, (i:N))
5 S^ = raise_dim(S, ((i * D) + j), (i:N,j:D))
6 C^ = raise_dim(C, j, (i:N,j:D))
7 for i in range(0, N):
8   min_sum!2[!I] = PHI(min_sum!1^[!I], min_sum!4[!I])
9   min_idx!2[!I] = PHI(min_idx!1^[!I], min_idx!4[!I])
10  sum!2 = [0,...,0]
11  sum!2^ = raise_dim(sum!2, (j:D))
12  d[!I,j] = SUB(S^[!I,j], C^[!I,j])
13  p[!I,j] = MUL(d[!I,j], d[!I,j])
14  for j in range(0, D):
15    sum!3[!I,j] = PHI(sum!2^[!I,j], sum!4[!I,j-1])
16    sum!4[!I,j] = ADD(sum!3[!I,j], p[!I,j])
17    sum!3^ = drop_dim(sum!3)
18    t[!I] = CMP(sum!3^[!I], min_sum!2[!I])
19    min_sum!3 = sum!3^
20    min_idx!3 = i
21    min_sum!4[!I] = MUX(t[!I], min_sum!3[!I], min_sum!2[!I])
22    min_idx!4[!I] = MUX(t[!I], min_idx!3[!I], min_idx!2[!I])
23 min_sum!2^ = drop_dim(min_sum!2)
24 min_idx!2^ = drop_dim(min_idx!2)
25 return (min_sum!2^, min_idx!2^)

```

When processing the outer loop two closures arise, one for $\text{min_sum!2}[!I] = \text{PHI}(\dots)$ and one for $\text{min_idx!2}[!I] = \text{PHI}(\dots)$. Since the two closures *do not* intersect, we have two distinct for-loops on i :

```

1 min_sum!1 = MAX_INT
2 min_sum!1^ = raise_dim(min_sum!1, (i:N))
3 min_idx!1 = 0
4 min_idx!1^ = raise_dim(min_idx!1, (i:N))

```

```

1161 5 S^ = raise_dim(S, ((i + D) + j), (i:N,j:D))
1162 6 C^ = raise_dim(C, j, (i:N,j:D))
1163 7
1164 8 sum!2 = [0,...,0]
1165 9 sum!2^ = raise_dim(sum!2, (j:D))
1166 10 d[l,j] = SUB(S^[l,j], C^[l,j])
1167 11 p[l,j] = MUL(d[l,j], d[l,j])
1168 12
1169 13 for j in range(0, D):
1170 14   sum!3[l,j] = PHI(sum!2^[l,j], sum!4[l,j-1])
1171 15   sum!4[l,j] = ADD(sum!3[l,j], p[l,j])
1172 16
1173 17 sum!3^ = drop_dim(sum!3)
1174 18 min_idx!3 = [0,1,2,...,N-1] // i.e., min_idx!3 = [i, (i:N)]
1175 19 min_sum!3 = sum!3^
1176 20
1177 21 for i in range(0, N):
1178 22   min_sum!2[i] = PHI(min_sum!1^[i], min_sum!4[i-1])
1179 23   t[i] = CMP(sum!3^[i], min_sum!2[i])
1180 24   min_sum!4[i] = MUX(t[i], min_sum!3[i], min_sum!2[i])
1181 25
1182 26 for i in range(0, N):
1183 27   min_idx!2[i] = PHI(min_idx!1^[i], min_idx!4[i-1])
1184 28   min_idx!4[i] = MUX(t[i], min_idx!3[i], min_idx!2[i])
1185 29
1186 30 min_sum!2^ = drop_dim(min_sum!2)
1187 31 min_idx!2^ = drop_dim(min_idx!2)
1188 32 return (min_sum!2^, min_idx!2^)
```

Phase 3 of Vectorization Algorithm. This phase removes redundant dimensionality. It starts by removing redundant dimensions in MOTION loops followed by removal of redundant drop dimension statements. It then does (extended) constant propagation to "bypass" raise statements, followed by copy propagation and dead code elimination.

The code becomes closer to what we started with:

```

1198 1 min_sum!1 = MAX_INT
1199 2 min_idx!1 = 0
1200 3 S^ = raise_dim(S, ((i + D) + j), (i:N,j:D))
1201 4 C^ = raise_dim(C, j, (i:N,j:D))
1202 5
1203 6 sum!2 = [0,...,0]
1204 7 d[l,j] = SUB(S^[l,j], C^[l,j])
1205 8 p[l,j] = MUL(d[l,j], d[l,j])
1206 9
1207 10 // j is redundant for sum!3 and sum!4
1208 11 for j in range(0, D):
1209 12   sum!3[l] = PHI(sum!2[l], sum!4[l])
1210 13   sum!4[l] = ADD(sum!3[l], p[l,j])
1211 14
1212 15 // drop_dim is redundant, removing
1213 16 // then copy propagation and dead code elimination
1214 17 min_idx!3 = [0,1,2,...,N-1] // i.e., min_idx!3 = [i, (i:N)]
1215 18
1216 19 // i is redundant for min_sum!2, min_sum!4 but not for t[i]
1217 20 for i in range(0, N):
```

```

21   min_sum!2 = PHI(min_sum!1, min_sum!4)
22   t[i] = CMP(sum!3[i], min_sum!2)
23   min_sum!4 = MUX(t[i], sum!3[i], min_sum!2)
24
25 // same, i is redundant for min_idx!2, min_idx!4
26 for i in range(0, N):
27   min_idx!2 = PHI(min_idx!1, min_idx!4)
28   min_idx!4 = MUX(t[i], min_idx!3[i], min_idx!2)
29
30 // drop_dim becomes redundant
31 return (min_sum!2, min_idx!2)
```

Phase 4 of Basic Vectorization. This phase adds SIMD operations:

```

1 min_sum!1 = MAX_INT
2 min_idx!1 = 0
3 S^ = raise_dim(S, ((i + D) + j), (i:N,j:D))
4 C^ = raise_dim(C, j, (i:N,j:D))
5
6 sum!2 = [0,...,0]
7 d[l,j] = SUB_SIMD(S^[l,j], C^[l,j])
8 p[l,j] = MUL_SIMD(d[l,j], d[l,j])
9
10 for j in range(0, D):
11   // l dim is a noop. sum is already a one-dimensional vector
12   sum!3[l] = PHI(sum!2[l], sum!4[l])
13   sum!4[l] = ADD_SIMD(sum!3[l], p[l,j])
14
15 min_idx!3 = [0,1,...,N-1]
16
17 for i in range(0, N):
18   min_sum!2 = PHI(min_sum!1, min_sum!4)
19   t[i] = CMP(sum!3[i], min_sum!2)
20   min_sum!4 = MUX(t[i], sum!3[i], min_sum!2)
21
22 for i in range(0, N):
23   min_idx!2 = PHI(min_idx!1, min_idx!4)
24   min_idx!4 = MUX(t[i], min_idx!3[i], min_idx!2)
25
26 return (min_sum!2, min_idx!2)
```

5.4 Correctness Argument

We build a correctness argument as follows. First, we define the MPC Source syntax. We then define the *linearization* of an MPC Source program as an *interpretation* over the syntax. The linearization is a *schedule* as defined in §3. We prove a theorem that states that the Basic vectorization algorithm preserves the def-use relations, or in other words, linearization of the vectorized MPC Source program gives rise to the exact same set of def-use pairs as linearization of the original program does. It follows easily that the schedule corresponding to the vectorized program computes the same result as the schedule corresponding to the original program.

MPC Source Syntax. Fig. 2 states the syntax and linearization semantics of MPC Source. Although notation is heavy, the linearization simply produces schedules as discussed in §2 and §3 the iterative MPC Source gives rise to what we called sequential

1277	$s ::= s_1; s_2$	$\gamma(s) = \gamma(s_1) ; \gamma(s_2)$	<i>sequence</i>	1335
1278	$ x[i, J, k] = \text{op_SIMD}(\gamma_1[i, J, k], \gamma_2[i, J, k])$	$\gamma(x[i, J, k] = \text{op_SIMD}(\gamma_1[i, J, k], \gamma_2[i, J, k])) =$	<i>operation</i>	1336
1279		$x[i, 0, k] = \gamma_1[i, 0, k] \text{ op } \gamma_2[i, 0, k] \ $		1337
1280		$x[i, 1, k] = \gamma_1[i, 1, k] \text{ op } \gamma_2[i, 1, k] \ \dots \ $		1338
1281		$x[i, J-1, k] = \gamma_1[i, J-1, k] \text{ op } \gamma_2[i, J-1, k]$		1339
1282	$ x[i, J, k] = \text{const}$	analogous	<i>constant</i>	1340
1283	$ x[i, J, k] = \text{PHI}(x_1[i, J, k], x_2[i, J, k-1])$		<i>pseudo PHI</i>	1341
1284	$ x[i, J, k] = \text{raise_dim}(x'[i], (J:J, k:K))$		<i>raise dimension(s)</i>	1342
1285	$ x[i, J] = \text{drop_dim}(x'[i, J, k], k)$		<i>drop dimension(s)</i>	1343
1286	$ \text{for } i \text{ in range}(I) : s$	$\gamma(\text{for } i \text{ in range}(I) : s) =$	<i>loop</i>	1344
1287		$\gamma(s)[0/i] ; \gamma(s)[1/i] ; \dots ; \gamma(s)[I-1/i]$		1345
1288				1346

Figure 2: MPC Source Syntax and Semantics. γ defines the semantics of MPC source which is a linearization of MPC Source. A SIMD operation parallelizes operations across the vectorized J dimension. \parallel denotes parallel execution, which is standard. γ of a for loop unrolls the loop. $;$ denotes sequential execution. Iterative MPC Source trivially extends to non-vectorized dimensions over the enclosing loops.

schedule where loops are unrolled and MPC Source with vectorized dimensions gives rise to what we called parallel schedule. For simplicity, we consider only scalars and read-only arrays, however, the treatment extends to write arrays as well (with our restriction on array writes to canonical writes). $x[i, J, k]$ denotes the value of scalar variable x at loop nest i, j, k . Upper case J denotes a vectorized dimension and lower case i, k denote iterative dimensions. There are semantic restrictions over the syntax: (1) x is treated as a 3-dimensional array and (2) $x[i, J, k]$ is enclosed into for-loops on non-vectorized dimensions i and k :

```

1305 1 for i in range(I):
1306 2   ...
1307 3   for k in range(K):
1308 4     ... x[i,j,k] ...

```

Linearization. Linearization is the concretization operation, which, as we mentioned earlier computes a schedule. The concretization function γ is defined as an interpretation of MPC Source syntax, as it is standard. The concretization of each one of the base statements is as follows:

$$\begin{aligned} \gamma(x[i, J, k] = \text{op_SIMD}(\gamma_1[i, J, k], \gamma_2[i, J, k])) = \\ x[i, 0, k] = \gamma_1[i, 0, k] \text{ op } \gamma_2[i, 0, k] \ || \\ x[i, 1, k] = \gamma_1[i, 1, k] \text{ op } \gamma_2[i, 1, k] \ || \dots \ || \\ x[i, I-1, k] = \gamma_1[i, I-1, k] \text{ op } \gamma_2[i, I-1, k] \end{aligned}$$

meaning that the vectorized dimension(s) are expanded into *parallel* statements. \parallel introduces SIMD (parallel) execution.

The concretization of the FOR statement is as follows:

$$\gamma(\text{FOR } 0 \leq i < I : s) = \gamma(s)[0/i] ; \gamma(s)[1/i] ; \dots \gamma(s)[I-1/i]$$

γ simply unrolls the loop substituting i with 0, 1, etc. Here $;$ denotes *sequential* execution.

As an example, consider the vectorized MPC Source excerpt from our running example. All variables are two dimensional arrays and the loop is vectorized in I but iterative in j :

```

1331 1 for j in range(0, D):
1332 2   sum!3[l,j] = PHI(sum!2[l,j], sum!4[l,j-1])
1333 3   sum!4[l,j] = ADD(sum!3[l,j], p[l,j])

```

Assuming $D = 2$ and $I = 2$ for simplicity, linearization produces the following schedule:

```

1355 1 sum!3[0,0] = PHI(sum!2[0,0], sum!4[0,-1]) ||
1356 2                                     sum!3[1,0] = PHI(sum!2[1,0], sum!4[1,-1])
1357 3 ;
1358 4 sum!4[0,0] = ADD(sum!3[0,0], p[0,0]) ||
1359 5                                     sum!4[1,0] = ADD(sum!3[1,0], p[1,0])
1360 6 ;
1361 7 sum!3[0,1] = PHI(sum!2[0,1], sum!4[0,0]) ||
1362 8                                     sum!3[1,1] = PHI(sum!2[1,1], sum!4[1,0])
1363 9 ;
1364 10 sum!4[0,1] = ADD(sum!3[0,1], p[0,1]) ||
1365 11 sum!4[1,1] = ADD(sum!3[1,1], p[1,1])

```

Note that by definition of the pseudo PHI function, $\text{PHI}(\text{sum}!2[\emptyset, \emptyset], \text{sum}!4[\emptyset, -1])$ evaluates to $\text{sum}!2[\emptyset, \emptyset]$ and therefore, the -1 index in the second argument does not matter.

Statements and def-uses over MPC Source. For each MPC Source program a we compute the def-use edges in the standard way: if base statement $s_0 \in a$ defines variable x , e.g., $x = \dots$, and base statement $s_1 \in a$ uses x , e.g., $\dots = \dots x$ and there is a path in the trivial CFG from s_0 to s_1 , then there is a def-use edge from s_0 to s_1 . We extend the dimensionality of a statement into $s_0[i, j, k]$ for statement s_0 enclosed in loop nest i, j, k .

Let a_0, a_1 be two MPC Source programs. Two statements, $s_0 \in a_0$ and $s_1 \in a_1$ are *same*, written $s_0 \equiv s_1$ if they are of the same operation and they operate on the same variables: same variable name and same dimensionality. Recall that dimensions in MPC Source are either iterative, lower case, or vectorized, upper case. Two statements are same even if one operates on an iterative dimension and the other one operates on a vectorized one, e.g., $s_0[i, j, k] \equiv s_1[I, j, K]$.

Statements and def-uses over linearized schedule. An *atomic* statement is a statement produced by linearization. We write $\underline{s_0}$ to denote statements in the concrete schedule as well as $\underline{s_0}[i, j, k]$ to denote fully instantiated values of i, j , and k , such as for example $\underline{s_0}[0, 1, 0]$. Clearly, the linearization of same statements produces the same set of atomic statements in the linearized schedule.

A def-use pair of atomic statements, denoted $(\underline{s}_0, \underline{s}_1)$ (indexing implicit), is defined in the standard way as well: \underline{s}_0 writes a location, say $x[\underline{i}, \underline{j}, \underline{k}]$, and \underline{s}_1 reads the same location.

Formal treatment. Property P defined below relates the linearized schedule of iterative MPC Source program a to the linearized schedule of the vectorized program a' . More precisely, a is the MPC Source program augmented with raise and drop dimension statements, i.e., Phase 1 without optimistic vectorization of all dimensions. a' is produced from a by Phase 2 of Basic Vectorization algorithm. We also write $\alpha(\underline{s})$ to map a statement back from $\gamma(a)$ to a , or from $\gamma(a')$ to a' .

DEFINITION 1. We say that $\gamma(a) \equiv \gamma(a')$ iff (1) atomic statement $\underline{s}[\underline{i}, \underline{j}, \underline{k}] \in \gamma(a)$ iff $\underline{s}[\underline{i}, \underline{j}, \underline{k}] \in \gamma(a')$ and (2) $(\underline{s}_0, \underline{s}_1) \in \gamma(a)$ iff $(\underline{s}_0, \underline{s}_1) \in \gamma(a')$ (indexing implicit).

The following lemma, which states that Basic vectorization preserves statements and def-use edges in the original MPC Source.

LEMMA 1. For each statement s in a , there is same statement s' in a' , and vice versa. For each def-use edge e in a , there is a same edge e' in a' , and vice versa.

PROOF. Proof sketch of Lemma 1. Phase 2 of Basic vectorization does not introduce any new statements in the code, it just vectorizes dimensions. Similarly, reordering of statements preserves exactly the def-use edges in the original MPC Source. \square

The main theorem is that property P holds for a and a' :

THEOREM 1. $\gamma(a) \equiv \gamma(a')$.

PROOF. Proof sketch of Theorem 1. The first condition of property P follows directly from Lemma 1. The proof of the second condition is by case-by-case analysis of the def-use edges in $\gamma(a)$.

To establish the second condition, which states that each def-use pair in $\gamma(a)$ is also in $\gamma(a')$ and there are no other pairs in $\gamma(a')$, first assume, without loss of generality, that there are no undefined uses in $\gamma(a)$. That is, each use is preceded by a definition, and inputs are special definitions.

From Lemma 1, it follows that the set of atomic statements in $\gamma(a)$ is exactly the same as in $\gamma(a')$. Therefore, we must show that for every $(\underline{s}_0, \underline{s}_1)$ in $\gamma(a)$, (1) \underline{s}_0 is scheduled before \underline{s}_1 in $\gamma(a')$ and (2) there is no other definition of the same variable scheduled between \underline{s}_0 and \underline{s}_1 .

Suppose there exists $(\underline{s}_0, \underline{s}_1)$ in $\gamma(a)$ such as either (1) or (2) is violated in $\gamma(a')$. We enumerate all cases and show contradiction.

First, consider a def-use pair $\underline{s}_0[\underline{i}, \underline{j}, \underline{k}] : x[\underline{i}, \underline{j}, \underline{k}] = \dots$ and $\underline{s}_1 : \dots = x$ (i.e., a read of same variable x) in the same iteration $\underline{i}, \underline{j}, \underline{k}$. Such an edge in the schedule $\gamma(a)$ could only result from MPC Source a that contains a same-level edge from $\alpha(\underline{s}_0) = s_0$ to $\alpha(\underline{s}_1) = s_1$. This edge is preserved by Basic vectorization in a' . Therefore, all possible $\gamma(a')$ will schedule \underline{s}_0 before \underline{s}_1 . To argue that (2) cannot be violated either, we again consider the possible cases and in all cases \underline{s}_0 will be made to read the value of x written at $\underline{s}_0[\underline{i}, \underline{j}, \underline{k}]$. Importantly, since $\gamma(a)$ is generated from SSA form, there is no other statement that writes the value of x , so all writes to x happen in previous or later iterations of $\underline{i}, \underline{j}, \underline{k}$. If the k dimension is not vectorized for either s_0 or s_1 , then clearly, there is no statement between the def

and the use. If it is vectorized for the definition s_0 but not for s_1 , then the algorithm ensures that the read references the appropriate value through indexing.

A similar case arises when \underline{s}_1 is in a subsequent iteration of the loop over k . This means that there is a backward edge from $\alpha(\underline{s}_0)$ to $\alpha(\underline{s}_1)$ in a and \underline{s}_1 is a PHI function. Analysis is similar as the backward edge in a' disallows scheduling of \underline{s}_1 before \underline{s}_0 . Since \underline{s}_1 occurs in the iteration that immediately follows the definition, \underline{s}_1 will read the correct value.

Second, consider a def-use pair $\underline{s}_0[\underline{i}, \underline{j}, \underline{k}] : x[\underline{i}, \underline{j}, \underline{k}] = \dots$ and $\underline{s}_1 : \dots = x$ in an enclosing iteration $\underline{i}, \underline{j}$. In this case, the value read at \underline{s}_1 is the value of x written at the last iteration $K-1$. Again, this edge is only possible if there is an inner-to-outer edge in MPC Source a and this edge is again preserved during vectorization. (*drop_dim* ensures that the read statement reads the last value written in x in both $\gamma(a)$ and $\gamma(a')$).

The last case is a def-use $\underline{s}_0[\underline{i}, \underline{j}] : x[\underline{i}, \underline{j}] = \dots$ and $\underline{s}_1 : \dots = x$ in an enclosed iteration $\underline{i}, \underline{j}, \underline{k}$. This again is possible only when there is an outer-to-inner edge in a . Reasoning is analogous to the previous cases and in this case *raise_dim* ensures that the value is replicated across all iterations of loop k . \square

It is a corollary of this theorem that the schedule a gives rise to (this is $\gamma(a)$), computes the exact same result as the schedule a' gives rise to (this is $\gamma(a')$):

COROLLARY 2. $\gamma(a)$ and $\gamma(a')$ produce same result, or more precisely, for every location $l[\underline{i}, \underline{j}, \underline{k}]$, $\gamma(a)$ and $\gamma(a')$ compute the same result.

PROOF. Proof sketch of Corollary 2. This can be established by induction over the length of def-use chains of computation in $\gamma(a)$. Assume that for all chains of length $\leq n$ all locations $l[\underline{i}, \underline{j}, \underline{k}]$ hold the same value in $\gamma(a)$ and $\gamma(a')$. A chain of length $n+1$ results from the execution of a statement $x[\underline{i}, \underline{j}, \underline{k}] = y[\underline{i}, \underline{j}, \underline{k}] \text{ op } z[\underline{i}, \underline{j}, \underline{k}]$. By property P , there is the same statement in $\gamma(a')$ and it is scheduled after the definitions of $y[\underline{i}, \underline{j}, \underline{k}]$ and $z[\underline{i}, \underline{j}, \underline{k}]$. By the inductive hypothesis $y[\underline{i}, \underline{j}, \underline{k}]$ and $z[\underline{i}, \underline{j}, \underline{k}]$ hold the same values in $\gamma(a')$ as in $\gamma(a)$. Therefore, locations $x[\underline{i}, \underline{j}, \underline{k}]$ hold the same value as well. We remark that due to the SSA form, each location $l[\underline{i}, \underline{j}, \underline{k}]$ is defined at most once. For clarity, we elide PHI nodes and raising and dropping dimensions; extending def-use reasoning is straight forward. \square

5.5 Extension with Array Writes

5.5.1 Removal of Infeasible Edges

Array writes limit vectorization as they sometimes introduce infeasible loop-carried dependencies. Consider the classical example from [3]:

```

1 for i in range(N):
2   A[i] = B[i] + 10;
3   B[i] = A[i] * D[i-1];
4   C[i] = A[i] * D[i-1];
5   D[i] = B[i] * C[i];

```

In Cytron's SSA this code (roughly) translates into

```

1509 1 for i in range(N):
1510 2   A_1 = PHI(A_0,A_2)
1511 3   B_1 = PHI(B_0,B_2)
1512 4   C_1 = PHI(C_0,C_2)
1513 5   D_1 = PHI(D_0,D_2)
1514 6   A_2 = update(A_1, i, B_1[i] + 10);
1515 7   B_2 = update(B_1, i, A_2[i] * D_1[i-1]);
1516 8   C_2 = update(C_1, i, A_2[i] * D_1[i-1]);
1517 9   D_2 = update(D_1, i, B_2[i] * C_2[i]);

```

There is a cycle around $B_1 = \text{PHI}(B_0, B_2)$ that includes statement $A_1 = \text{update}(A_0, i, B_1[i] + 10)$ and that statement won't be vectorized even though in fact there is no loop-carried dependency from the write of $B_1[i]$ at 7 to the read of $\dots = B_1[i]$ at 6.

The following algorithm removes certain infeasible loop-carried dependencies that are due to array writes. Consider a loop with index $0 \leq j < J$ nested at i, j, k . Here i represents the enclosing loops of j and k represents the enclosed loops in j .

```

1527 for each array A written in loop j do
1528   { including enclosed loops in j }
1529   dep = False
1530   for each pair def:  $A_m[f(i, j, k)] = \dots$ , and use:  $\dots = A_n[f'(i, j, k)]$  in loop j do
1531     if  $\exists \underline{i}, \underline{j}', \underline{k}, \underline{k}'$ , s.t.  $0 \leq \underline{i} < I$ ,  $0 \leq \underline{j}' < J$ ,  $0 \leq \underline{k}, \underline{k}' < K$ ,  $\underline{j} < \underline{j}'$ ,
1532       and  $f(\underline{i}, \underline{j}, \underline{k}) = f'(\underline{i}, \underline{j}', \underline{k}')$  then
1533         dep = True
1534       end if
1535     end for
1536   end for
1537   if dep == False then
1538     remove back edge into A's  $\phi$ -node in loop j.
1539   end if
1540 end for

```

Consider a loop j enclosed in some fixed \underline{i} . Only if an update (definition) $A_m[f(i, j, k)] = \dots$ at some iteration \underline{j} references the same array element as a use $\dots = A_n[f'(i, j, k)]$ at some later iteration \underline{j}' , we may have a loop-carried dependence for A due to this def-use pair. (In contrast, Cytron's algorithm inserts a loop-carried dependency every time there is an array update.) The algorithm above examines all def-use pairs in loop j , including defs and uses in nested loops, searching for values $\underline{i}, \underline{j}', \underline{k}, \underline{k}'$ that satisfy $f(\underline{i}, \underline{j}, \underline{k}) = f'(\underline{i}, \underline{j}', \underline{k}')$. If such values exist for some def-use pair, then there is a potential loop-carried dependence on A; otherwise there is not and we can remove the spurious backward edge thus "freeing up" statements for vectorization.

More precisely, we use Z3 [17] to check satisfiability of the formula

$$(0 \leq \underline{i} < I) \wedge (0 \leq \underline{j}' < J) \wedge (0 \leq \underline{k}, \underline{k}' < K) \wedge (\underline{j} < \underline{j}') \wedge f(\underline{i}, \underline{j}, \underline{k}) = f'(\underline{i}, \underline{j}', \underline{k}')$$

Formulas f and f' are simple as loop nests are typically of depth 2-3. Therefore, Z3 completes the process instantly.

Consider the earlier example. There is a single loop, i . Clearly, there is no pair \underline{i} and \underline{i}' , where $\underline{i} < \underline{i}'$ that make $\underline{i} = \underline{i}'$ due to the def-use pairs of A 6-7 and 6-8. Therefore, we remove the back edge from 6 to the phi-node 2. Analogously, we remove the back edges from 7 to 3 and from 8 to 4. However, there are many values $\underline{i} < \underline{i}'$ that make $\underline{i} = \underline{i}' - 1$ and the back edge from 9 to 5 remains (def-use pairs for D). As a result of removing these spurious edges,

Vectorization will find that statement 6 is vectorizable. Statements 7, 8 and 9 will correctly appear in the FOR loop.

Note however, that this step renders some array phi-nodes target-less. We handle target-less phi-nodes with a minor extension of Vectorization (Phase 2). First, we merge closures that update the same array. This simplifies handling of array ϕ -nodes: if each closure is turned into a separate loop each loop will need to have its own array phi-node to account for the update and this would complicate the analysis. Second, we add the target-less node of array A back to the closure that updates A — the intuition is, even if there is no loop-carried dependence from writes to reads on A, A is written and the write (i.e., update) cannot be vectorized; therefore, the updated array has to carry to the next iteration of the loop. Third, in cases when the phi-node remains target-less, i.e., cases when the array write can be vectorized, we have to properly remove the phi-node replacing uses of the left-hand side of the phi-node with its arguments.

5.5.2 Restricting Array Writes

For now, we restrict array updates to *canonical updates*. Assume (for simplicity) a two-dimensional array $A[I, J]$. A canonical update is the following:

```

1589 1 for i in range(I):
1590 2   for j in range(J):
1591 3     ...
1592 4     A[i,j] = ...
1593 5     ...

```

Here i and j are the outermost loops. The update $A[i, j]$ can be nested into an inner loop and there may be multiple updates, i.e., writes to $A[i, j]$. However, update such as $A[i-1, j] = \dots$ or $A[i-1, j-1] = \dots$, etc., is not allowed. Additionally, while there could be several different loops that perform canonical updates, they must be of the same dimensionality, i.e., an update of higher or lower dimension, e.g., $A[i, j, k] = \dots$ is not allowed. We compute the *canonical dimensionality* of each write array by examining the array writes in the original program and rejecting programs that violate the canonical write restriction. This restriction simplifies reasoning in this early stage of the compiler; we will look to relax the restriction in future work.

The reason for this restriction is that it allows us to easily extend correctness reasoning from the basic case in the previous section. The restriction ensures that the array shape does not change and raise dimension and drop dimension can be applied in the same way as in the basic case.

Reads through an arbitrary formula, such as $A[i-1]$ for example, are allowed; currently, the projection function returns dummy values if the read formula is out of bounds; we assume the programmer ensures that the program still computes correct output in this case.

5.5.3 Changes to Basic Vectorization

In addition to the changes for the handling of target-less phi-nodes, Basic Vectorization has to handle def-use edges $X \rightarrow Y$ where X defines and Y uses an array variable. The definition can be an update $A_2 = \text{update}(A_1, i, \dots)$, a pseudo ϕ -node $A_2 = \text{PHI}(A_0, A_1)$, etc.. Note that ϕ -nodes for arrays have no subscript operations the way there are subscript operations in analysis-introduced arrays representing scalars. While there are variations, the most intuitive

implementation will perform Basic Vectorization Phase 1 as is, inserting *raise_dim* and *drop_dim* in the same way. However, the implementation of raise dimension and drop dimension will be adapted because the dimension cannot be raised or dropped to a dimension lower than the canonical one. Consider a def-use edge $X \rightarrow Y$ for an array A.

- (1) same-level $X \rightarrow Y$. Do nothing, propagate the array, which happens to be of the right dimension.
- (2) inner-to-outer $X \rightarrow Y$ triggers the addition of *drop_dim*. However, the dimensionality cannot be dropped below the canonical dimensionality of the array. E.g., if the dimensionality of the loop enclosure X is already at the canonical one, then *drop_dim* has no effect.
- (3) outer-to-inner $X \rightarrow Y$ triggers *raise_dim*. Again, if the dimensionality of the loop enclosure of Y is smaller or same as the canonical dimensionality of the array, then it has no effect, otherwise, if dimensionality is greater than the canonical dimensionality, *raise_dim(...)* (at X) is the same as in Basic Vectorization.
- (4) "mixed" $X \rightarrow Y$. We assume that the mixed edge is transformed into an inner-to-outer followed by outer-to-inner edge before we perform vectorization, just as with Basic vectorization.

If the use of the array is a read $A[f(i, j, k)]$ different than a canonical read $A[i, j, k]$, then we need to add a reshape operation as all arrays are $A[i, j, k]$. It can be added after *raise_dim*/*drop_dim* or incorporated in these operations. The bulk of the change is in Phase 2 of Vectorization as outlined earlier.

5.5.4 Examples with Array Writes

Example 1. First, the canonical dimensionality of all A,B,C and D is 1. After Phase 1 of Vectorization the Aiken's array write example will be (roughly) as follows:

```

1 for i in range(N):
2   A_1 = PHI(A_0,A_2)
3   B_1 = PHI(B_0,B_2)
4   C_1 = PHI(C_0,C_2)
5   D_1 = PHI(D_0,D_2)
6   A_2 = update(A_1, l, B_1[l] + 10);
7   B_2 = update(B_1, l, A_2[l] * D_1[l-1]);
8   C_2 = update(C_1, l, A_2[l] * D_1[l-1]);
9   D_2 = update(D_1, l, B_2[l] * C_2[l]);
    
```

Note that since all def-uses are same-level (i.e., reads and writes of the array elements) no raise dimension or drop dimension happens.

Phase 2 computes the closure of 5; $cl = \{5, 7, 8, 9\}$ while 6 is vectorizable. Recall that 2,3, and 4 are target-less phi-nodes. Since the closure cl includes updates to B and C, the corresponding phi-nodes are added back to the closure and the def-use edges are added back to the target-less nodes. The uses of A_1 and B_1 in the vectorized statement turn into uses of A_0 and B_0 respectively; this is done for all original target-less phi-node. (But note that A_0 is irrelevant; the update writes into array A_2 in parallel.) Finally, the target-less phi-node for A is discarded.

```

1 A_2 = update(A_0, l, ADD_SIMD(B_0[l],10));
2   equiv. to A_2[l] = ADD_SIMD(B_0[l],10)
3 for i in range(N): // MOTION loop
    
```

```

4   B_1 = PHI(B_0,B_2)
5   C_1 = PHI(C_0,C_2)
6   D_1 = PHI(D_0,D_2)
7   B_2 = update(B_1, i, A_2[i] * D_1[i-1]);
8     equiv. to B_2 = B_1; B_2[i] = A_2[i] * D_1[i-1];
9   C_2 = update(C_1, i, A_2[i] * D_1[i-1]);
10  D_2 = update(D_1, i, B_2[i] * C_2[i]);
    
```

Example 2. Now consider the MPC Source of Histogram:

```

1 for i in range(0, num_bins):
2   res1 = PHI(res, res2)
3   for j in range(0, N):
4     res2 = PHI(res1, res3)
5     tmp1 = (A[j] == i)
6     tmp2 = (res2[i] + B[j])
7     tmp3 = MUX(tmp1, res2[i], tmp2)
8     res3 = Update(res2, i, tmp3)
9 return res1
    
```

The canonical dimensionality of res is 1. Also, the phi-node $res1 = PHI(res, res2)$ is a target-less phi-node (the implication being that the inner for loop can be vectorized across i). After Phase 1, Vectorization produces the following code (statements are implicitly vectorized along i and j). In a vectorized update statement, we can ignore the incoming array, res2 in this case. The update writes (in parallel) all locations of the 2-dimensional array, in this case it sets up each $res3[i,j] = tmp3[i,j]$.

```

1 A1 = raise_dim(A, j, ((i:num_bins),(j:N)))
2 B1 = raise_dim(B, j, ((i:num_bins),(j:N)))
3 l = raise_dim(i, ((i:num_bins),(j:N)))
4 for i in range(0, num_bins):
5   res1 = PHI(res, res2^)# target-less phi-node
6   res1^ = raise_dim(res1, (j:N))
7   for j in range(0, N):
8     res2 = PHI(res1^, res3)
9     tmp1 = (A1 == l)
10    tmp2 = (res2 + B1)
11    tmp3 = MUX(tmp1, res2, tmp2)
12    res3 = Update(res2, (l,j), tmp3)
13    res2^ = drop_dim(res2)
14 res1'' = drop_dim(res1)
15 return res1''
    
```

Processing the inner loop in Phase 2 vectorizes $tmp1 = (A1 == l)$ along the j dimension but leaves the rest of the statements in a MOTION loop. Processing the outer loop is interesting. This is because the PHI node is a target-less node, and therefore, there are no closures! Several things happen. (1) Everything can be vectorized along the i dimension. (2) We remove the target-less PHI node, however, we must update uses of res1 appropriately: the use at *raise_dim* goes to the first argument of the PHI function and the use at *drop_dim* goes to the second argument.

```

1 A1 = raise_dim(A, j, ((i:num_bins),(j:N)))
2 B1 = raise_dim(B, j, ((i:num_bins),(j:N)))
3 l1 = raise_dim(i, ((i:num_bins),(j:N)))
4
5 tmp1[l,j] = (A1[l,j] == l1[l,j])
    
```

```

1741 6
1742 7 res1^ = raise_dim(res, (j:N)) // replacing res1 with res, 1st arg
1743 8 for j in range(0, N):
1744 9   res2 = PHI(res1^, res3)
1745 10  tmp2[1,j] = (res2[1,j] + B1[1,j])
1746 11  tmp3[1,j] = MUX(tmp1[1,j], res2[1,j], tmp2[1,j])
1747 12  res3 = Update(res2, (1,j), tmp3)
1748 13  equiv. to res3 = res2; res3[1,j] = tmp3[1,j]
1749 14  res2^ = drop_dim(res2)
1750 15  res1 = drop_dim(res2^) // replacing with res2^, 2nd arg. NOOP
1751 16  return res1

```

6 COMPILER BACK END

MOTION code generation requires that variables are marked as `plain` or `shared` following the type system in §4.1. We require that all inputs are marked as either `shared` or `plaintext`, however, we infer qualifiers for the rest of the variables. §6.1 briefly describes the taint analysis and §6.2 describes MOTION code generation.

6.1 Taint Analysis

The taint analysis works on MPC Source, which lacks if-then-else control flow. This significantly simplifies treatment as there is no need to handle conditionals and implicit flow. Specifically, the compiler uses the following rules, which are standard in positive-negative qualifier systems (here `shared` is the positive qualifier and `plain` is the negative one):

- (1) Loop counters are always `plain`.
- (2) If any variable on the right-hand side `rhs` of an assignment is `shared`, then the assigned variable `lhs` is `shared` following subtyping rule `rhs <: lhs`.
- (3) Any variables that cannot be determined as `shared` via the above rules are `plain`.

In the below snippet `sum!2` and `sum!3` form a dependency cycle and there is no `shared` value that flows to either one. They are inferred as `plaintext`.

```

1777 1 plaintext_array = [0, 1, 2, ...]
1778 2 sum!1 = 0
1779 3 for i in range(0, N):
1780 4   sum!2 = PHI(sum!1, sum!3)
1781 5   sum!3 = sum!2 + plaintext_array[i]

```

When converting to MOTION code, any `plaintext` value used in the right-hand side of a `shared` assignment is converted to a `shared` value for that expression.

6.2 From (Optimized) MPC Source to MOTION

MOTION supports FOR loops and SIMD operations, so translation from MPC source to MOTION C++ code is relatively straightforward.

Variable declarations: Our generated C++ uses the following variable-naming scheme: `shared` variables are named the same as in the MPC Source with the `!` replaced with an underscore (e.g. `sum!2` would be translated to `sum_2`). `Plaintext` variables follow the same naming convention as `shared` variables but are prefixed with `_MPC_PLAINTEXT_`. The `shared` representation of constants are

named `_MPC_CONSTANT_` followed by the literal constant (e.g. the `shared` constant `0` would be named `_MPC_CONSTANT_0`).

The generated MOTION code begins with the declaration of all variables used in the function, including loop counters. If a variable is a vectorized array, it is initialized to a correctly-sized array of empty MOTION shares. Additionally, each `plaintext` variable and parameter has a `shared` counterpart declared. Next, all constant values which are used as part of `shared` expressions are initialized as a `shared` input from party 0. Finally, `plaintext` parameters are converted used as `shared` inputs from party 0 to initialize their `shared` counterparts.

Code generation: Once the function preamble is complete, the MPC Source is translated into C++ one statement at a time. The linear structure of MPC Source enables this approach to translation. If there is no vectorization present in a statement, translation to C++ is straightforward: outside of MUX statements and array updates, non-vectorized assignments, expressions, and returns directly translate into their C++ equivalents. Non-vectorized MUX statements are converted to MOTION's MUX member function on the condition variable. Array updates are translated into two C++ assignments: one to update the value in the original array and one to assign the new array as shown in Listing 2.

MPC FOR loops are converted to C++ FOR loops which iterate the loop counter over the specified range. Pseudo PHI nodes are broken into two components: the "FALSE" branch which assigns the initial value of the PHI node and the "TRUE" branch which assigns the PHI node's back-edge. The assignment of the "FALSE" branch occurs right before the PHI node's enclosing loop. As these assignments may rely on the loop counter, the loop counter is initialized before these statements. Inside of the PHI node's enclosing loop, a C++ `if` statement is inserted to only assign the true branch of the PHI node after the first iteration. Listing 3 illustrates this translation.

Vectorization and SIMD operations: Vectorization is handled with utility functions to manage accessing and updating slices of arrays. All SIMD values are stored in non-vectorized form as 1-dimensional `std::vectors` in row-major order. Whenever a SIMD value is used in an expression, the utility function `vectorized_access()` takes the multi-dimensional representation of a SIMD value, along with the size of each dimension and the requested slice's indices, and converts that slice to a MOTION SIMD value. Because MOTION supports SIMD operations using the same C++ operators as non-SIMD operations, we do not need to perform any other transformations to the expression. Therefore, once vectorized accesses are inserted the translation of an expression containing SIMD values is identical to that of expressions without SIMD values.

Similarly, the `vectorized_assign()` function assigns a (potentially SIMD) value to a slice of a vectorized array. This operation cannot be done with a simple subscript as SIMD assignments will update a range of values in the underlying array representation.

Updating SIMD arrays is also implemented differently from updating non-vectorized arrays. Instead of separating the array update from the assignment of the new array, these steps are combined with the `vectorized_update()` utility function. This function operates identically to `vectorized_assign()`, however it additionally returns the array after the assignment occurs. This value is then used for the assignment to the new variable. Listing 4 illustrates

1857 1858 1859 1860 1861 1862 1863 1864 1865 1866 1867 1868 1869 1870 1871 1872 1873 1874 1875 1876 1877 1878 1879 1880 1881 1882 1883 1884 1885 1886 1887 1888 1889 1890 1891 1892 1893 1894 1895 1896 1897 1898 1899 1900 1901 1902 1903 1904 1905 1906 1907 1908 1909 1910 1911 1912 1913 1914	<pre> 1 A[i] = val </pre> <p style="text-align: center;">IMP Source</p>	<pre> 1 A!2 = update(A!1, i, val) </pre> <p style="text-align: center;">MPC Source</p>	<pre> 1 A_1[i] = val; 2 A_2 = A_1; </pre> <p style="text-align: center;">MOTION Code</p>	1915 1916 1917 1918 1919 1920 1921 1922 1923 1924 1925 1926 1927 1928 1929 1930 1931 1932 1933 1934 1935 1936 1937 1938 1939 1940 1941 1942 1943 1944 1945 1946 1947 1948 1949 1950 1951 1952 1953 1954 1955 1956 1957 1958 1959 1960 1961 1962 1963 1964 1965 1966 1967 1968 1969 1970 1971 1972
--	---	--	--	--

Table 2: MOTION Translation: Array Updates

1864 1865 1866 1867 1868 1869 1870 1871 1872 1873 1874 1875 1876 1877 1878 1879 1880 1881 1882 1883 1884 1885 1886 1887 1888 1889 1890 1891 1892 1893 1894 1895 1896 1897 1898 1899 1900 1901 1902 1903 1904 1905 1906 1907 1908 1909 1910 1911 1912 1913 1914	<pre> 1 for i in range(N): 2 tmp = PHI(arr[i], val!0) 3 ... </pre> <p style="text-align: center;">MPC Source</p>	<pre> 1 _MPC_PLAINTEXT_i = 0; 2 tmp = arr[_MPC_PLAINTEXT_i]; 3 for (; _MPC_PLAINTEXT_i < _MPC_PLAINTEXT_N; _MPC_PLAINTEXT_i++) { 4 if (_MPC_PLAINTEXT_i != 0) { 5 tmp = val_0; 6 } 7 ... 8 } </pre> <p style="text-align: center;">MOTION Code</p>	1915 1916 1917 1918 1919 1920 1921 1922 1923 1924 1925 1926 1927 1928 1929 1930 1931 1932 1933 1934 1935 1936 1937 1938 1939 1940 1941 1942 1943 1944 1945 1946 1947 1948 1949 1950 1951 1952 1953 1954 1955 1956 1957 1958 1959 1960 1961 1962 1963 1964 1965 1966 1967 1968 1969 1970 1971 1972
--	--	---	--

Table 3: MOTION Translation: FOR loop with Phi nodes

1876 1877 1878 1879 1880 1881 1882 1883 1884 1885 1886 1887 1888 1889 1890 1891 1892 1893 1894 1895 1896 1897 1898 1899 1900 1901 1902 1903 1904 1905 1906 1907 1908 1909 1910 1911 1912 1913 1914	<pre> 1 sum!4[l] = ADD_SIMD(sum!3[l], p[l, j]) </pre> <p style="text-align: center;">MPC Source</p>	<pre> 1 vectorized_assign(sum_4, {_MPC_PLAINTEXT_N}, {true}, {}), 2 vectorized_access(sum_3, {_MPC_PLAINTEXT_N}, {true}, {}) + 3 vectorized_access(p, {_MPC_PLAINTEXT_N, _MPC_PLAINTEXT_D}, {true, false}, 4 {_MPC_PLAINTEXT_j}); </pre> <p style="text-align: center;">MOTION Code</p>	1915 1916 1917 1918 1919 1920 1921 1922 1923 1924 1925 1926 1927 1928 1929 1930 1931 1932 1933 1934 1935 1936 1937 1938 1939 1940 1941 1942 1943 1944 1945 1946 1947 1948 1949 1950 1951 1952 1953 1954 1955 1956 1957 1958 1959 1960 1961 1962 1963 1964 1965 1966 1967 1968 1969 1970 1971 1972
--	---	---	--

Table 4: MOTION Translation: Assignment to SIMD value

vectorized_assign() and vectorized_update() on the Biometric example.

Reshaping and raising dimensions: Raising the dimensions of a scalar or array uses the lift() utility function which takes a lambda for the raised expression and the dimensions of the output. This function is also used for the scalar expansion of values which have been lifted out of FOR loops as described in §5.2. This function evaluates the expression for each permutation of indices along the dimensions and returns the resulting array in row-major order. The lambda accepts an array of integers representing the index along each of the dimensions being raised, and the translation of the expression which is being raised replaces each of the dimension index variables with the relevant subscript of this array. There is also a special case of the lift() function which occurs when we are raising an array. In this case, instead of concatenating the array for each index, we extend the array along all dimensions being raised which are not present in the array already. For example, when raising an array with dimensions $N \times M$ to an array with dimensions $N \times M \times D$, the input array will simply be extended along the D dimension: $A'[n, m, d] = A[n, m]$ for every d . If the input array is already correctly sized it will be returned as-is.

Dropping dimensions use the drop_dim() and drop_dim_monoreturn() utility functions. They function identically but the latter returns a scalar for the case when the final dimension of an array is dropped. These functions take the non-vectorized representation of an array, along with the

dimensions of that array, and return the array with the final dimension dropped.

Upcasting from plaintext to shared: Currently, our compiler only supports the Bmr and BooleanGMW protocols as MOTION does not implement all operations for other protocols. MOTION does not support publicly-known constants for these protocols, so all conversions from plaintext values to shares are performed by providing the plaintext value as a shared input from party 0. Due to this limitation, our translation to MOTION code attempts to minimize the number of conversions from a plaintext value. This is accomplished by creating a shared copy of each plaintext variable and updating that copy in lock-step with the plaintext variable. Since variables are often initialized to a common constant value (e.g. 0), this approach decreases the number of input gates by only creating a shared input for each initialization constant. Loop counters must still be converted to a shared value on each iteration that they are used, however we only generate this conversion when necessary, i.e., when the counter flows to a shared computation. This is to prevent unnecessary increase in the number of input gates when loop counters are only used as plaintext.

Due to the SSA translation phase as well as the conversions to and from SIMD values which our utility functions perform, our generated vectorized MOTION code often includes multiple copies of arrays and scalar values. These copies do not incur a runtime cost as the arrays simply hold pointers to the underlying shares, so no new shares or gates are created as a result of this copying.

```
1 raise_dim(i + j, (i:N, j:M))
```

MPC Source

```
1 lift(std::function([&](const std::vector<std::uint32_t> &idxs) {return idxs[0] + idxs[1]}),
2     {_MPC_PLAINTEXT_N, _MPX_PLAINTEXT_M})
```

MOTION Code

Table 5: MOTION Translation: Raising dimensions

Cost in MPC programs is dominated by shares and computation on shares.

7 EXPERIMENTAL RESULTS

7.1 Experiment Setup

We have tested our framework with several benchmarks. For the multiparty computation (MPC), we restricted our evaluation to 2 party computation (2PC) setting because it requires fewer computing resources. We stress that there is no such inherent restriction in our framework. We use hardware resources provided by CloudLab[20] and consider two network settings, namely the Local Area Network (LAN) and the Wide Area Network (WAN). In the LAN setting, we use c6525-25g machines connected via a 10Gbps link with sub-millisecond latency. These machines are equipped with 16-core AMD 7302P 3.0GHz processors and 128GB of RAM. This setting reflects typical LAN use-case considering that 10Gbps LAN is increasingly common in business networks and is now available even in some home networks. In the WAN setting we again used a c6525-25g machine (located in Utah, US) for the first party and a c220g1 machine (located in Wisconsin, US) for the second. The c220g1 machine is equipped with two Intel E5-2630 8-core 2.40GHz processors and 128GB of RAM. We measured the connection bandwidth between these machines to be 560Mbps and average round trip time (RTT) to be 38ms. At the time of this writing, all major internet providers in the US offer 1Gbps connections to home consumers, therefore this setting reasonably reflects the typical WAN use case.

We run all experiments 5 times and report average values of various metrics. Note that standard deviation, shown as error bar on top of the histogram bars in the graphs, in all observations was at most 4.5% of the mean value, therefore the accuracy of the results is not effected by the relatively fewer runs.

7.2 Benchmarks

In the following, we say *both* for an experiment in which we execute both non-vectorized and vectorized protocols and *vec* for the vectorized only experiment. For each problem, we run benchmark experiment with increasing sizes of input e.g. we ran biometric matching with database size N of $\{2, 4, 8, \dots, 4096\}$; At some input size 2^k (e.g. $N = 2^8$ in biometric matching), the non-vectorized version runs out of memory. From this point on, we run only *vec* experiment up to input size of $2^{2(k-1)}$ e.g. in case of biometric matching, we run the *vec* experiment up to database size $N = 2^{12}$. The *vec* experiment completes without issues (e.g. running out of memory) for input size of $2^{2(k-1)}$ for nearly all benchmarks. While the numbers for all runs are not shown for space reasons, the run times are largely consistent, e.g. if non-vectorized experiment fails

at some input size 2^k , and *vec* experiment takes X seconds to complete for the input size 2^{k-1} , then for an input size $2^{k+\ell-1}$, it takes roughly $\ell \cdot X$ seconds to complete.

We use the following benchmarks in our evaluation:

- (1) *Biometric Matching*: Server has a database S of N records, each record's dimension is D . Client submits a query C , client and server compute the closest record to C in an MPC. We use $N=128$ for *both* and $N=4096$ for *vec*. D is fixed at 4.
- (2) *Convex Hull*: Given a polygon of N vertices (split between Alice and Bob), convex hull is computed in an MPC. It is adapted from [22]. We use $N=32$ for *both* experiment and $N=256$ for *vec* experiment.
- (3) *Count 102*: Alice has a string of length N of symbols, Bob has a regular expression of the form $1(0^*)2$, together they compute number of substrings that match the regular expression. It is adapted from [22]. We use $N=1024$ for *both* and $N=4096$ for *vec*.
- (4) *Count 10*: Same as *Count 102* except now the regular expression is of the form $1(0^+)$. Parameters are same as above.
- (5) *Cryptonets Max Pooling*: Given a matrix of rows \times cols elements that are split between Alice and Bob, they compute the max pooling subroutine of the cryptonet benchmark[19]. We use rows=64, cols=64 for *both* experiment.
- (6) *Database Join*: given two databases with A and B containing 2-element records, compute cross join. We use $A=B=32$ for *both* and $A=B=64$ for *vec*.
- (7) *Database Variance* given a database of len records, compute variance. We use $len=512$ for *both* and $len=4096$ for *vec*.
- (8) *Histogram*: Given N 5-star ratings, compute their histogram, taken from [25, 22]. We use $N=512$ for *both* and $N=4096$ for *vec*.
- (9) *Inner Product*: given two vectors, each of N elements, compute their inner product. We use $N=512$ for *both* and $N=4096$ for *vec*.
- (10) *k-means Iteration*: performs the iteration subroutine of k-means database clustering operation [26, 37]. Here $len1$ is the size of input data, and $len2$ is the number of clusters. We use $len1=32$, $len2=5$ for *both* and $len1=256$, $len2=8$ for *vec*.
- (11) *Longest 102*: Similar to *Count 102* except that it computes the largest substring matching the regular expression. We use same parameters as *Count 102*, adapted from [22].
- (12) *Max Distance b/w Symbols* Alice has a string of N symbols and Bob has some symbol θ . The MPC computes the maximum distance between θ s in the string. We adapted it from [22]. We use $N=1024$ for *both* and $N=2048$ for *vec*.
- (13) *Minimal Points* Given a set of N points (split between Alice and Bob), a set of minimal points is computed i.e. there is no

other point that has both a lower x and y coordinate, adapted from [22]. We use $N=32$ for *both* and $N=64$ for *vec*.

- (14) *MNIST ReLU* given an input of outer \times inner elements, executes the MNIST ReLU subroutine. We use $\text{inner}=512$ for *both* and $\text{inner}=2048$ for *vec*. outer is fixed at 16.
- (15) *Private Set Intersection (PSI)* Alice holds set S_1 with size SA , Bob holds set S_2 with size SB , together they compute intersection of their sets. We use $SA=SB=128$ for *both* and $SA=SB=1024$ for *vec*.

7.3 Results and Analysis

A detailed summary of the effects of vectorization on various benchmarks is presented in Table 6. We show circuit evaluation times in Fig. 3. In terms of amenability to vectorization, we divide benchmarks into 3 categories: 1) *High*: these include convex hull, cryptonets max pooling, minimal points and private set intersection. These benchmarks are highly parallelizable and see 70x to 25x speedup in BMR, and 55x to 30x in GMW protocol. 2) *Medium*: these include biometric matching, DB Variance, histogram, inner product, k-means iteration and MNIST ReLU. These benchmarks have non-parallelizable phases e.g. the summing phase of inner product and biometric matching. Still, most computation is parallelizable and it results in speedup from 25x to 5x in BMR, and 25x to 2x in GMW protocol. 3) *Low*: these include the database join and the regular expression benchmarks (count 102, count 10, longest 102 and max distance between symbols). Very few operations in these programs are parallelizable, thus the speedup is lower. We see a speedup from 2x to 1.1x in BMR. In GMW, database join, count 102 and count 10s see speedup from 1.3x to 1.1x. However, longest102 and max distance between symbols suffer a slowdown of 0.5x. There is some opportunity for vectorization in these benchmarks according to our analytical model, particularly, there is a large EQ operation that is vectorized, although a large portion of the loop cannot be vectorized. We observe that transformation to vectorized code increases multiplicative depth and, the negative effect of increased depth is more noticeable in a round-based protocol like GMW. The cause of the increase is not clear – we conjecture that MOTION performs optimizations over the non-vectorized loop body that decreases depth; also, EQ is relatively inexpensive in Boolean GMW and BMR compared to ADD and MUL, which also de-emphasizes the benefit of vectorization. We propose a simple heuristic although we do leave all the benchmarks in the table: if the transformation increases circuit depth beyond some threshold (e.g. more than 10% of the original circuit), we can reject the transformation. Note that in some settings it may still be desirable to vectorize e.g. in data constrained environments. As shown in Fig. 8, vectorization results in reduced communication (fewer bits are transferred).

This reduction is a result of more efficient data-packing at both 1) the application level (i.e. the MPC framework level) and 2) at the network level. The MPC framework needs to store/send metadata with each primitive/message so that it can correctly decoded/consumed later. For example, a gate needs an identifier *gid*, a gate type *gtype*, incoming wire identifiers, etc. Say $\text{size}(gmeta)$, bits are needed to store/send metadata for a single gate. Using one vectorized/SIMD gate instead of $(N + 1)$ non-vectorized gates saves $N \cdot \text{size}(gmeta)$ bits in memory/communication. Similarly, at the network level,

each message needs a header h that contains routing and decoding information while the packet is in transit. Say, one (non-vectorized) interactive gate induces a payload p . This means, $\text{size}(h) + \text{size}(p)$ bits are sent to network for each (non-vectorized) interactive gate. Evaluation of N such gates translates to $N \cdot (\text{size}(h) + \text{size}(p))$ bits of communication. On the other hand, a vectorized gate that replaces these N gates is much cheaper, requiring $\text{size}(h) + N \cdot \text{size}(p)$ bits of communication.

Concretely, let us consider an MPC framework implemented on Transport Control Protocol (TCP) over Internet Protocol (IP). This is the most common communication stack in today's applications. Moreover, for the sake of communication size comparison, the only difference between UDP and TCP is the smaller header size of 8 bytes² in UDP compared to the at least 20 bytes³ in TCP. Both protocols are typically implemented over Internet Protocol and the header size of an IPv4 packet is 20 bytes⁴. In the Arithmetic GMW protocol, multiplication operation (MUL) is typically implemented using Beaver's triples [?]. This means that, in the online phase, all parties need to send 2ℓ bits to each other. If $\ell = 32$ bits, then, TCP payload is $2\ell = 64$ bits or 8 bytes. Considering that the Maximum Transmission Unit (MTU) is typically 1500 bytes, a TCP message may have payload of up to $(1500 - 20 - 20) = 1460$ bytes (the exact value for this size is decided via the Maximum Segment Size (MSS) during TCP stack initialization, Specification maximum is 65,496 bytes). Meaning, a vectorized gate replacing $1460/8 \approx 180$ non-vectorized gates could be sent in a single 1500 byte message rather than $182 \cdot (20 + 20 + 8) = 8,640$ bytes required otherwise. Similar reasoning applies to interactive gates in boolean GMW and, while exact improvement depends on the implementation details, packing data reduces both the memory and communication footprint regardless of the underlying MPC framework (as long as it supports vectorized gates).

In the case of BMR, the entire circuit can be packed as one payload and sent using a few TCP packets. Therefore under-utilization of network's payload-capacity is not an issue. At the application (MPC framework) level however, packing inefficiency can still be a problem. For example, MOTION uses 64 bits (8 bytes) for gate identifiers. A vectorized gate that replaces 128 non-vectorized gates, requires only one gate identifier i.e. 8 bytes. On the other hand, storing 128 non-vectorized gates requires 1,024 bytes (an additional 1,016 bytes) to store 128 gate identifiers. Thus, vectorization reduces the size of the circuit. This, in turn, reduces payload for the network and means that fewer TCP packets need to be sent, thereby saving on TCP/IP metadata that would have been needed for additional packets.

In summary, vectorization provides better packing for interactive protocols like GMW compared to a constant round protocol like BMR and, indeed, this is the trend we see in Fig. 8. Fig. 10 confirms that the number of gates are fewer in vectorized circuits. Consequently, in the highly vectorizable benchmarks, circuit generation time (see Fig. 9) for vectorized circuits is a fraction of non-vectorized circuit because a much smaller circuit is being generated. For the

²https://en.wikipedia.org/wiki/User_Datagram_Protocol#UDP_datagram_structure

³https://en.wikipedia.org/wiki/Transmission_Control_Protocol#TCP_segment_structure

⁴<https://en.wikipedia.org/wiki/IPv4#Header>

Table 6: Vectorized vs Non-Vectorized Comparison, times in seconds (in LAN setting where applicable), Communication in MiB, Numbers in 1000s, values rounded to nearest integer, benchmark names ending in V are vectorized.

Benchmark	GMW						BMR					
	Online	Setup	# Gates	Circ Gen	# Msgs	Comm.	Online	Setup	# Gates	Circ Gen	# Msgs	Comm.
Biometric Matching	146	16	1,784	119	1,413	140	89	263	1,595	139	2,716	312
Biometric Matching (V)	12	4	34	2	28	14	2	13	30	4	61	130
Convex Hull	48	6	551	40	516	51	28	72	494	39	695	80
Convex Hull (V)	0	1	2	0	1	4	0	2	1	1	2	32
Count 102	79	6	418	35	525	52	15	62	269	33	785	92
Count 102 (V)	71	5	316	24	332	34	11	30	167	16	304	59
Count 10s	79	6	419	35	525	52	14	62	270	33	785	92
Count 10s (V)	71	4	316	24	332	34	11	29	167	16	304	59
Cryptonets (Max Pooling)	50	11	688	46	554	55	36	89	608	51	898	110
Cryptonets (Max Pooling) (V)	1	1	7	1	2	5	2	4	7	2	12	49
Database Join	70	8	433	48	790	80	19	229	458	119	3,518	427
Database Join (V)	54	6	320	35	575	61	16	112	320	57	1,457	285
Database Variance	166	18	2,009	135	1,639	163	95	269	1,708	145	2,795	320
Database Variance (V)	37	6	321	24	334	43	10	30	170	13	178	141
Histogram	94	10	862	68	979	97	27	94	491	51	1,132	135
Histogram (V)	33	5	166	16	164	23	7	17	92	13	154	68
Inner Product	127	15	1,675	108	1,308	130	83	250	1,526	134	2,623	301
Inner Product (V)	16	5	158	12	165	25	6	18	83	7	86	127
k-means	108	12	1,333	88	1,090	108	63	185	1,141	99	1,958	225
k-means (V)	6	3	47	4	43	12	2	11	32	4	54	95
Longest 102	93	7	650	52	713	71	26	93	475	49	1,091	128
Longest 102 (V)	169	6	544	41	519	53	25	60	369	33	605	95
Max. Dist. b/w Symbols	71	8	572	43	576	57	24	69	397	38	748	89
Max. Dist. b/w Symbols (V)	166	7	538	39	512	51	24	57	363	32	589	78
Minimal Points	35	5	458	31	369	37	24	46	401	26	347	40
Minimal Points (V)	0	1	1	0	1	3	0	1	1	0	1	16
MNIST ReLU	132	31	1,843	126	1,483	152	98	247	1,630	135	2,401	298
MNIST ReLU (V)	3	3	25	3	9	17	5	11	25	5	33	136
Private Set Intersection	95	9	558	59	1,049	104	22	186	591	96	2,639	302
Private Set Intersection (V)	1	2	1	2	1	8	1	8	2	4	2	122

benchmarks that are not amenable to vectorization, circuit generation is still shorter, although not substantially so. Online time and setup time are presented in Fig. 11, and Fig. 12 respectively.

Let us look more closely at the Biometric Matching benchmark in Fig. 4, Fig. 5, and Fig. 6. For input size beyond $N=128$ the memory usage exceeds available memory and prevents circuit generation. Consequently, non-vectorized bars are missing beyond this threshold in the graphs. Notice that vectorization improves all metrics. In circuit evaluation (see Fig. 4), BMR sees higher speedup (23x faster) compared to GMW (10x faster), while GMW sees faster circuit generation time at 45x lower (see Fig. 6) compared to BMR's which is

35x lower. Communication size reduction (see Fig. 5) is higher for GMW (10x less) compared to BMR (2.5x less).

Since our vectorization framework is network agnostic, it produces the same circuit for both LAN and WAN. This means that the number of gates and communication size remain the same. Moreover, time for circuit generation, which is a local operation, also does not change. Setup and Online times, however, increase due to lower bandwidth and higher latency of the WAN. Indeed, this is what we observe in Fig. 7.

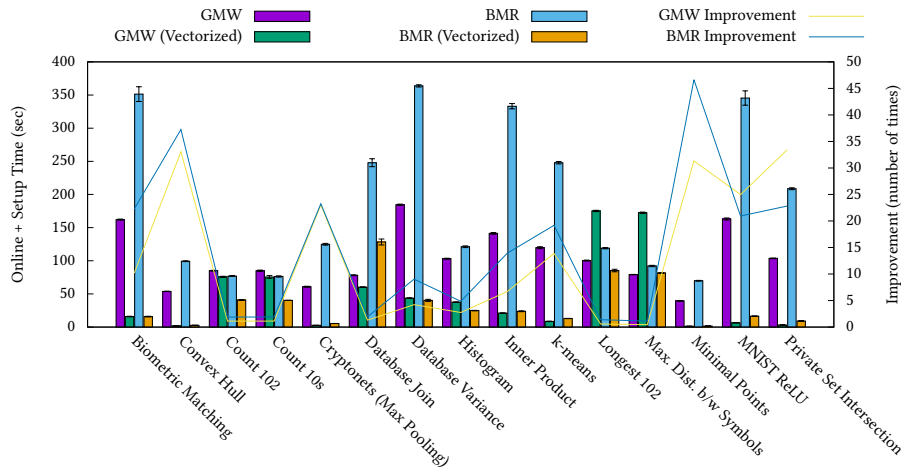


Figure 3: Circuit Evaluation Time (Setup and Online Phase) in Seconds, LAN Setting. The Error-bars are Standard Deviation.

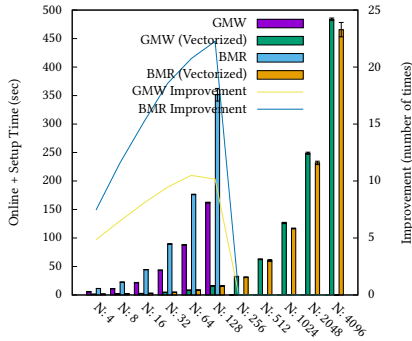


Figure 4: Biometric Matching Circuit Evaluation Time, x-axis lists database size

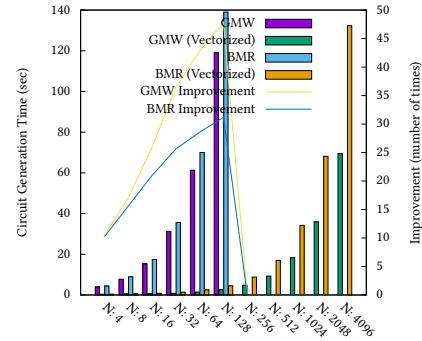


Figure 6: Biometric Matching Circuit Generation Time, x-axis lists database size

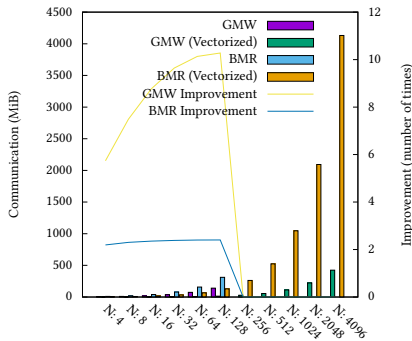


Figure 5: Biometric Matching Communication Size, x-axis lists database size

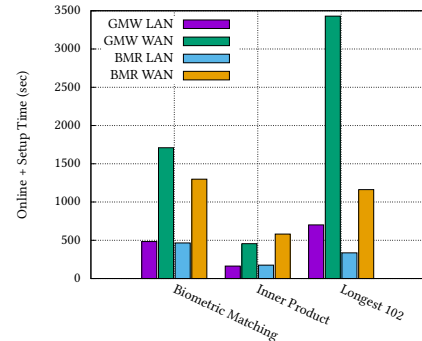


Figure 7: LAN vs. WAN: Circuit Evaluation Time Comparison

7.4 Comparison with MOTION-native Inner Product

Finally, we compared our automatically generated routine for Inner Product with the manually SIMD-ified MOTION-native routine in the distribution. We were surprised that we were an order of

magnitude slower in Boolean GMW as our circuit ran a significantly larger number of communication rounds. Upon investigation, it turns out that the vectorized multiplication are the same, however, our addition loop incurs significant cost (ADD is non-local and expensive in Boolean GMW). The MOTION-native loop runs

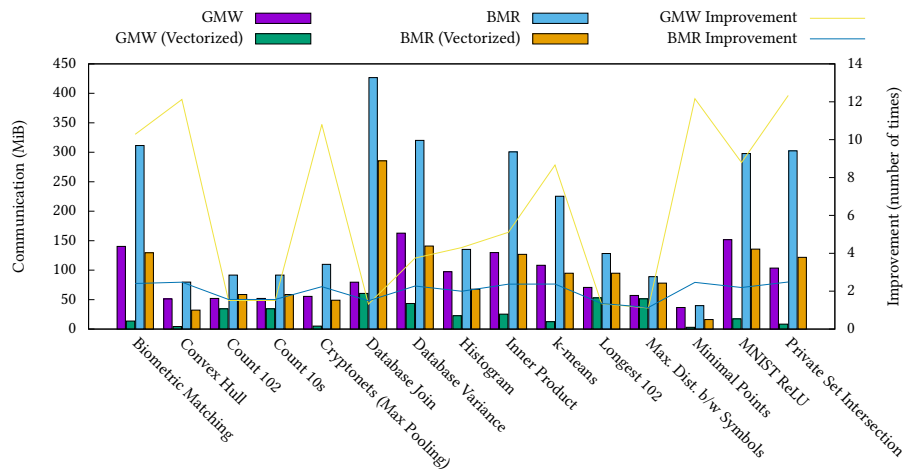


Figure 8: Communication Size of Benchmarks

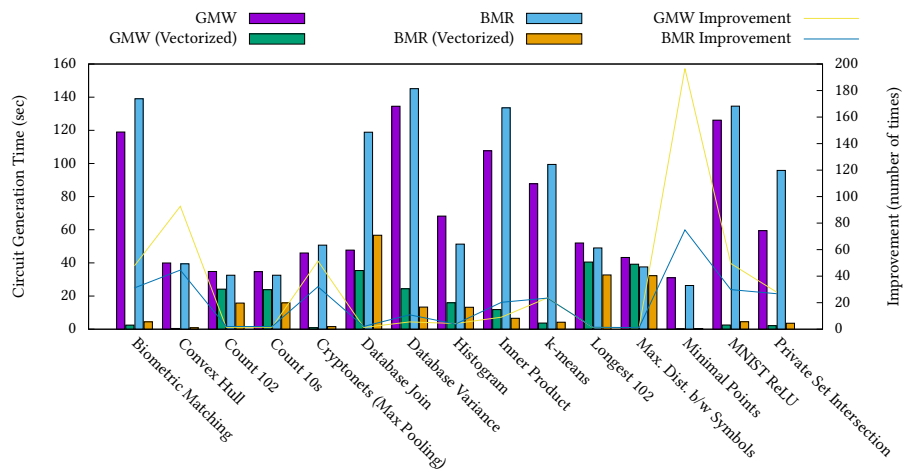


Figure 9: Circuit Generation Time of Benchmarks

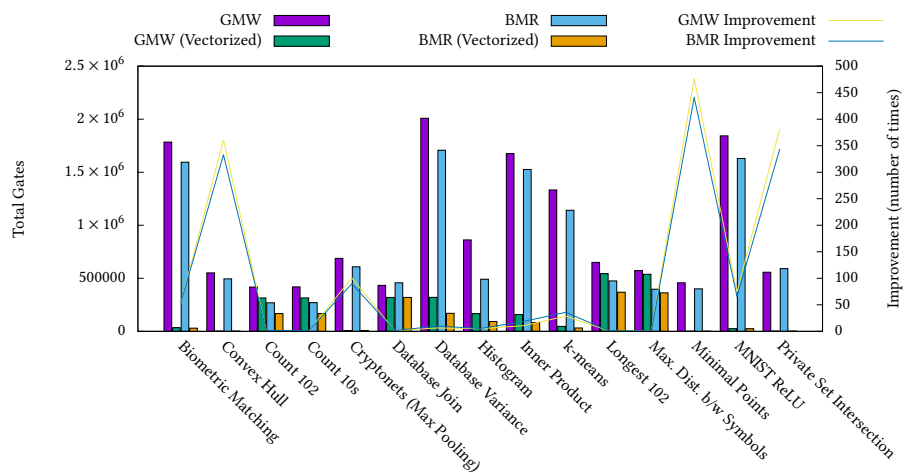


Figure 10: Number of Gates of Benchmarks

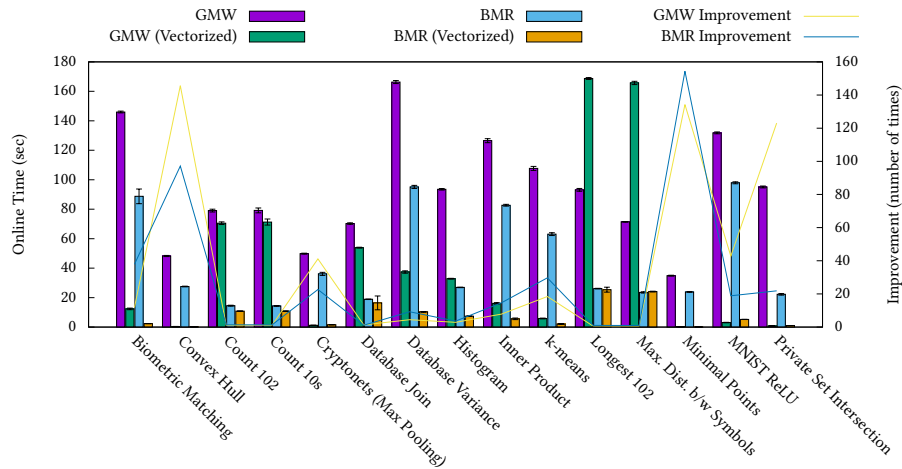


Figure 11: Online Time of Benchmarks

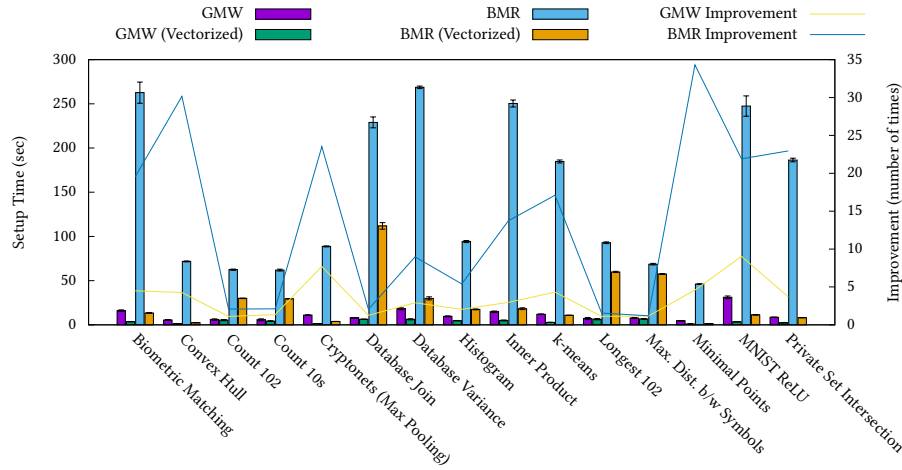


Figure 12: Setup Time of Benchmarks

```
1 result += mult_unsimplified[i];
```

while our loop generates and runs

```
1 result[i] = result[i-1] + mult_unsimplified[i];
```

Recall that the scalar expansion is an artifact of our vectorization. We rewrote the accumulation (manually, for testing purposes) and that lead to the same running time.

MOTION’s compiler performs analysis that informs circuit generation and the example illustrates the power of the analysis. In the above example, MOTION does the standard divide-and-conquer accumulation in $O(\log(N))$ rounds. Recall that the next phase of Vectorization, in §5, which we have not implemented yet, gets rid of redundant dimensions and generates

```
1 if (i != 0) {
2   result_prev = result;
3 }
4 result = result_prev + mult_unsimplified[i];
```

. And while it is unrealistic to expect that MOTION’s static analysis will detect the associative accumulation in the scalar expansion code, it is realistic to expect that it will in the above code. It still might lead to confusion in the static analysis as analysis on the AST is difficult. Our investigation showed that not only MOTION does not optimize

```
1 if (i!=0) {
2   result_prev = result;
3 }
4 result = result_prev + mult_unsimplified[i];
```

it does not optimize the much simpler accumulation:

```
1 result = result + mult_unsimplified[i];
```

We conjecture that MPC Source, a straight-forward representation, will not only enable detection of general associative loops, but also allow for program synthesis to increase opportunities

for divide-and-conquer parallelization [22]; as the problem is non-trivial, particularly the interaction of divide-and-conquer with vectorization and mixing, we leave it for future work.

8 RELATED WORK

MPC languages and compilers. Languages and compilers for secure computation have seen significant attention and advances in recent years. The early MPC compilers Fairplay [6], and Sharemind [10] were followed by PICCO [41], Obliv-C [40], TinyGarble [36], Wystiria [34], and others. A new generation of MPC compilers includes SPDZ/SCALE-MAMBA/MP-SPDZ [28] and the ABY/HyCC/MOTION [18, 14, 12] frameworks. These two families are the state-of-the-art and are actively developed. Another recent development is Viaduct, a functional language and compiler that supports a range of secure computation frameworks, including MPC and ZKP. Hastings et al. present a review of compiler frameworks [24].

While each of these languages and compilers brings in new ideas and advances, none addresses the problem of “circuit independent” intermediate representation and optimization. We envision a classical compiler structure: (1) a Wysteria, Viaduct, Obliv-C, or IMP Source front end, including rich type systems and AST-level semantic analysis, compile into the MPC Source IR, (2) MPC Source-level optimizations take place, followed by (3) back-end compilers into circuits. Our focus is at the intermediate level.

Many works focus on the implementation of MPC protocols exposing an API to the programmer. For example, the ABY/MOTION line of compiler frameworks provides a library of MPC primitives; the programmer writes MPC programs in C++ on top of the library. These back ends implement different protocols and allow for mixing, but notably, they leave it to the programmer to assign different protocols to different parts of the computation and perform share conversion accordingly. In addition, MOTION provides SIMD primitives, which allows for efficient execution of MPC operations, but again, using SIMD primitives is the responsibility of the programmer. There is interest in frameworks for automatic mixing, e.g., [14, 25, 21].

Other works, e.g., Obliv-C [40], Wystiria [34] and Viaduct [1] focus on higher-level language design, particularly information-flow systems that restrict flow between secure and insecure parts of the program. OblivM [29] has similar goals to ours — it optimizes oblivious (i.e., secure) source code into an efficient garbled circuit. Our work complements OblivM, in the sense that while OblivM relies on programmer annotations such as map-reduce constructs, we automatically detect opportunities for optimization at an intermediate level of representation. We have implemented vectorization, however, divide-and-conquer is a natural next step. Ozdemir et al. [33] develop CirC, an IR with backends into zero-knowledge proof primitives as well as SMT primitives. Our work focuses on MPC, which CirC does not support yet. MPC Source is a higher level of representation than CirC, e.g., it does not unroll loops, thus presenting a smaller-size input program suitable for vectorization and scalable program analysis. Overall, we view our work and CirC as orthogonal — we are focusing on analysis and optimization while CirC focuses on infrastructure.

HyCC [14] is a compiler from C Source into ABY circuits. It does source-to-source compilation with the key goal to decompose the program into modules and then assign protocols to modules. In contrast, we focus on MPC Source-level optimizations, specifically vectorization, although we envision future optimizations as well. We formalize MPC Source and reasoning about transformations, which we conjecture is more tractable than reasoning over the higher-level AST. On the other hand, HyCC does inter-procedural optimizations, while our analysis is intra-procedural. We will explore context-sensitive inter-procedural analysis over MPC Source in future work. HyCC, similarly to Buscher [13] uses an of-the-shelf source-to-source polyhedral compiler⁵ to perform vectorization at the level of source code. The disadvantage of using an of-the-shelf source-to-source compiler is that it solves a more general problem than what MPC presents and may forgo opportunities for optimization — concretely, it is well-known that vectorization and polyhedral compilation do not work well with conditionals [8, 27]. In contrast, we consider vectorization at the level of MPC Source which linearizes conditionals; we are able to handle programs with interleaved if-statements and for-loops and achieve significant speedup.

Classical HPC compilers. Automatic vectorization is a longstanding problem in high-performance computing (HPC). There are thousands of works in this area reflecting over 40 years of research. We presented a vectorization algorithm for MPC Source, essentially extending classical loop vectorization [4]. In HPC vectorization, conditional control flow presents a challenge — one cannot estimate the cost of a schedule or vectorize branches in a straightforward manner — in contrast to MPC Source vectorization. We view Karrenberg’s work on Whole function vectorization [27] as most closely related to ours — it linearizes the program and vectorizes both branches of a conditional applying masking to avoid execution of the branch-not-taken code, and selection (similar to MUX) to select the correct value based on the result of the condition at runtime. The problem is that masking and selection, or more generally, handling control predicates [8, 27], can lead to *slowdown*.

We argue that vectorization over linear MPC Source is a different problem, one that warrants a new look, while drawing from results in HPC. Since both branches of the conditional and the multiplexer *always* execute, not only can we apply aggressive vectorization on linear code, but (perhaps more importantly) we can also build analytical models that accurately predict execution time. These models in turn would drive optimizations such as vectorization, protocol mixing, and others. Vectorization meshes in with those additional optimizations in non-trivial ways.

Furthermore, extensions of classical loop vectorization with array writes, arbitrary indexing, including non-affine indexing, and interaction with SSA are non-trivial and present novel challenges and opportunities for contribution. Polyhedral parallelization [8] considers a higher-level source (typically AST) representation, while our work takes advantage of linear MPC Source and SSA form. The work by Karrenberg [27] is rare in that space, in the sense that it considers vectorization over SSA form, which has similarities to MPC Source. We consider different array representation, notion of

⁵We believe HyCC uses Par4All (<https://github.com/Par4All/par4all>), however, does not appear to be included with the publicly available distribution of HyCC.

dependence, and reasoning about dependence, which we conjecture is more suitable for MPC Source.

9 CONCLUSION AND FUTURE WORK

We presented a formalization of the MPC Source intermediate language followed by a specific back-end optimization at the level of MPC Source: novel SIMD-vectorization. We demonstrated that vectorization has significant impact on performance. We are excited about the opportunities for future work — integration with protocol mixing, divide-and-conquer reasoning and parallelization, as well as inter-procedural context-sensitive analysis at the level of MPC Source will improve MPC programmability and efficiency.

REFERENCES

[1] Cosku Acay, Rolph Recto, Joshua Gancher, Andrew C. Myers, and Elaine Shi. Viaduct: an extensible, optimizing compiler for secure distributed programs. In Stephen N. Freund and Eran Yahav, editors, *PLDI '21: 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation, Virtual Event, Canada, June 20-25, 2021*, pages 740–755. ACM, 2021.

[2] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison-Wesley series in computer science / World student series edition. Addison-Wesley, 1986.

[3] Alexander Aiken and Alexandru Nicolau. Optimal loop parallelization. In Richard L. Wexelblat, editor, *Proceedings of the ACM SIGPLAN'88 Conference on Programming Language Design and Implementation (PLDI), Atlanta, Georgia, USA, June 22-24, 1988*, pages 308–317. ACM, 1988.

[4] Randy Allen and Ken Kennedy. Automatic translation of fortran programs to vector form. *ACM Trans. Program. Lang. Syst.*, 9(4):491–542, 1987.

[5] Toshinori Araki, Assi Barak, Jun Furukawa, Marcel Keller, Yehuda Lindell, Kazuma Ohara, and Hikaru Tsuchida. Generalizing the SPDZ compiler for other protocols. In David Lie, Mohammad Mannan, Michael Backes, and XiaoFeng Wang, editors, *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS 2018, Toronto, ON, Canada, October 15-19, 2018*, pages 880–895. ACM, 2018.

[6] Assaf Ben-David, Noam Nisan, and Benny Pinkas. Fairplaymp: a system for secure multi-party computation. In *Proceedings of the 2008 ACM Conference on Computer and Communications Security, CCS 2008, Alexandria, Virginia, USA, October 27-31, 2008*, pages 257–266, 2008.

[7] Michael Ben-Or, Shafi Goldwasser, and Avi Wigderson. Completeness theorems for non-cryptographic fault-tolerant distributed computation (extended abstract). In Janos Simon, editor, *Proceedings of the 20th Annual ACM Symposium on Theory of Computing, May 2-4, 1988, Chicago, Illinois, USA*, pages 1–10. ACM, 1988.

[8] Mohamed-Walid Benabderrahmane, Louis-Noël Pouchet, Albert Cohen, and Cédric Bastoul. The polyhedral model is more widely applicable than you think. In Rajiv Gupta, editor, *Compiler Construction, 19th International Conference, CC 2010, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2010, Paphos, Cyprus, March 20-28, 2010. Proceedings*, volume 6011 of *Lecture Notes in Computer Science*, pages 283–303. Springer, 2010.

[9] Marina Blanton and Paolo Gasti. Secure and efficient protocols for iris and fingerprint identification. In *Computer Security - ESORICS 2011 - 16th European Symposium on Research in Computer Security, Leuven, Belgium, September 12-14, 2011. Proceedings*, pages 190–209, 2011.

[10] Dan Bogdanov, Sven Laur, and Jan Willemson. Sharemind: A framework for fast privacy-preserving computations. In *Computer Security - ESORICS 2008, 13th European Symposium on Research in Computer Security, Málaga, Spain, October 6-8, 2008. Proceedings*, pages 192–206, 2008.

[11] Peter Bogetoft, Dan Lund Christensen, Ivan Damgård, Martin Geisler, Thomas P. Jakobsen, Mikkel Kroigaard, Janus Dam Nielsen, Jesper Buus Nielsen, Kurt Nielsen, Jakob Pagter, Michael I. Schwartzbach, and Tomas Toft. Secure multiparty computation goes live. In Roger Dingledine and Philippe Golle, editors, *Financial Cryptography and Data Security, 13th International Conference, FC 2009, Accra Beach, Barbados, February 23-26, 2009. Revised Selected Papers*, volume 5628 of *Lecture Notes in Computer Science*, pages 325–343. Springer, 2009.

[12] Lennart Braun, Daniel Demmler, Thomas Schneider, and Oleksandr Tkachenko. Motion? a framework for mixed-protocol multi-party computation. *ACM Trans. Privacy and Security*, 25(2):1–35, May 2022.

[13] Niklas Büscher. *Compilation for More Practical Secure Multi-Party Computation*. PhD thesis, Darmstadt University of Technology, Germany, 2018.

[14] Niklas Büscher, Daniel Demmler, Stefan Katzenbeisser, David Kretzmer, and Thomas Schneider. Hycc: Compilation of hybrid protocols for practical secure computation. In David Lie, Mohammad Mannan, Michael Backes, and XiaoFeng Wang, editors, *Proceedings of the 2018 ACM SIGSAC Conference on Computer and*

Communications Security, CCS 2018, Toronto, ON, Canada, October 15-19, 2018, pages 847–861. ACM, 2018.

[15] David Chaum, Claude Crépeau, and Ivan Damgård. Multiparty unconditionally secure protocols (extended abstract). In Janos Simon, editor, *Proceedings of the 20th Annual ACM Symposium on Theory of Computing, May 2-4, 1988, Chicago, Illinois, USA*, pages 11–19. ACM, 1988.

[16] Ron Cytron, Jeanne Ferrante, Barry K. Rosen, Mark N. Wegman, and F. Kenneth Zadeck. Efficiently computing static single assignment form and the control dependence graph. *ACM Trans. Program. Lang. Syst.*, 13(4):451?–490, 1991.

[17] Leonardo Mendonça de Moura and Nikolaj S. Björner. Z3: an efficient SMT solver. In C. R. Ramkrishnan and Jakob Rehof, editors, *Tools and Algorithms for the Construction and Analysis of Systems, 14th International Conference, TACAS 2008, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2008, Budapest, Hungary, March 29-April 6, 2008. Proceedings*, volume 4963 of *Lecture Notes in Computer Science*, pages 337–340. Springer, 2008.

[18] Daniel Demmler, Thomas Schneider, and Michael Zohner. ABY - A framework for efficient mixed-protocol secure two-party computation. In *22nd Annual Network and Distributed System Security Symposium, NDSS 2015, San Diego, California, USA, February 8-11, 2015*.

[19] Nathan Dowlin, Ran Gilad-Bachrach, Kim Laine, Kristin Lauter, Michael Naehrig, and John Wernsing. Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, ICML'16*, page 201–210. JMLR.org, 2016.

[20] Dmitry Duplyakin, Robert Ricci, Aleksander Maricq, Gary Wong, Jonathon Duerig, Eric Eide, Leigh Stoller, Mike Hibler, David Johnson, Kirk Webb, Aditya Akella, Kuangching Wang, Glenn Ricart, Larry Landweber, Chip Elliott, Michael Zink, Emmanuel Cecchet, Snigdhaswin Kar, and Prabodh Mishra. The design and operation of CloudLab. In *Proceedings of the USENIX Annual Technical Conference (ATC)*, pages 1–14, July 2019.

[21] Vivian Fang, Lloyd Brown, William Lin, Wenting Zheng, Aurojit Panda, and Raluca Ada Popa. Costco: An automatic cost modeling framework for secure multi-party computation. Cryptology ePrint Archive, Report 2022/332, 2022. <https://ia.cr/2022/332>.

[22] Azadeh Farzan and Victor Nicolet. Phased synthesis of divide and conquer programs. In Stephen N. Freund and Eran Yahav, editors, *PLDI '21: 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation, Virtual Event, Canada, June 20-25, 2021*, pages 974–986. ACM, 2021.

[23] Oded Goldreich, Silvio Micali, and Avi Wigderson. How to play any mental game or A completeness theorem for protocols with honest majority. In *Proceedings of the 19th Annual ACM Symposium on Theory of Computing, 1987, New York, New York, USA*, pages 218–229, 1987.

[24] Marcella Hastings, Brett Hemenway, Daniel Noble, and Steve Zdancewic. Sok: General purpose compilers for secure multi-party computation. In *2019 IEEE Symposium on Security and Privacy, SP 2019, San Francisco, CA, USA, May 19-23, 2019*, pages 1220–1237. IEEE, 2019.

[25] Muhammad Ishaq, Ana L. Milanova, and Vassilis Zikas. Efficient MPC via program analysis: A framework for efficient optimal mixing. In Lorenzo Cavallaro, Johannes Kinder, XiaoFeng Wang, and Jonathan Katz, editors, *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security, CCS 2019, London, UK, November 11-15, 2019*, pages 1539–1556. ACM, 2019.

[26] Geetha Jagannathan and Rebecca N. Wright. Privacy-preserving distributed k-means clustering over arbitrarily partitioned data. In *Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining, KDD '05*, page 593–599, New York, NY, USA, 2005. Association for Computing Machinery.

[27] Ralf Karrenberg. *Automatic SIMD Vectorization of SSA-based Control Flow Graphs*. Springer, 2015.

[28] Marcel Keller. MP-SPDZ: A versatile framework for multi-party computation. In Jay Ligatti, Xinming Ou, Jonathan Katz, and Giovanni Vigna, editors, *CCS '20: 2020 ACM SIGSAC Conference on Computer and Communications Security, Virtual Event, USA, November 9-13, 2020*, pages 1575–1590. ACM, 2020.

[29] Chang Liu, Xiao Shaun Wang, Kartik Nayak, Yan Huang, and Elaine Shi. Oblivm: A programming framework for secure computation. In *2015 IEEE Symposium on Security and Privacy, SP 2015, San Jose, CA, USA, May 17-21, 2015*, pages 359–376. IEEE Computer Society, 2015.

[30] Payman Mohassel and Peter Rindal. ABY3: A mixed protocol framework for machine learning. In *Proc. ACM Conf. on Comp. Comm. Sec.*, pages 35?–52, New York, NY, USA, 2018. ACM.

[31] Payman Mohassel and Yupeng Zhang. Secureml: A system for scalable privacy-preserving machine learning. In *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*, pages 19–38. IEEE Computer Society, 2017.

[32] Tobias Nipkow and Gerwin Klein. *Concrete Semantics: With Isabelle/HOL*. Springer, Heidelberg, Germany, 2014.

[33] Alex Ozdemir, Fraser Brown, and Riad S. Wahby. Unifying compilers for snarks, smt, and more. *IACR Cryptol. ePrint Arch.*, page 1586, 2020.

2901	[34] Aseem Rastogi, Matthew A. Hammer, and Michael Hicks. Wysteria: A programming language for generic, mixed-mode multiparty computations. In <i>2014 IEEE Symposium on Security and Privacy, SP 2014, Berkeley, CA, USA, May 18-21, 2014</i> , pages 655–670, 2014.	2959
2902		2960
2903		2961
2904	[35] Michael L. Scott. <i>Programming Language Pragmatics (3. ed.)</i> . Academic Press, 2009.	2962
2905		2963
2906	[36] Ebrahim M. Songhori, Siam U. Hussain, Ahmad-Reza Sadeghi, Thomas Schneider, and Farinaz Koushanfar. Tinygarble: Highly compressed and scalable sequential garbled circuits. In <i>2015 IEEE Symposium on Security and Privacy, SP 2015, San Jose, CA, USA, May 17-21, 2015</i> , pages 411–428. IEEE Computer Society, 2015.	2964
2907		2965
2908	[37] Jaideep Vaidya and Chris Clifton. Privacy-preserving k -means clustering over vertically partitioned data. In <i>Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '03</i> , page 206–215, New York, NY, USA, 2003. Association for Computing Machinery.	2966
2909		2967
2910		2968
2911	[38] Vijay V. Vazirani. <i>Approximation Algorithms</i> . Springer, Heidelberg, Germany, 2010.	2969
2912		2970
2913	[39] Andrew C. Yao. Protocols for secure computations. In <i>Proceedings of the 23rd Annual Symposium on Foundations of Computer Science, SFCS '82</i> , pages 160–164, Washington, DC, USA, 1982. IEEE Computer Society.	2971
2914		2972
2915	[40] Samee Zahur and David Evans. Obliv-c: A language for extensible data-oblivious computation. Cryptology ePrint Archive, Report 2015/1153, 2015. https://ia.cr/2015/1153 .	2973
2916		2974
2917	[41] Yihua Zhang, Aaron Steele, and Marina Blanton. PICCO: a general-purpose compiler for private distributed computation. In <i>2013 ACM SIGSAC Conference on Computer and Communications Security, CCS'13, Berlin, Germany, November 4-8, 2013</i> , pages 813–826, 2013.	2975
2918		2976
2919		2977
2920		2978
2921		2979
2922		2980
2923		2981
2924		2982
2925		2983
2926		2984
2927		2985
2928		2986
2929		2987
2930		2988
2931		2989
2932		2990
2933		2991
2934		2992
2935		2993
2936		2994
2937		2995
2938		2996
2939		2997
2940		2998
2941		2999
2942		3000
2943		3001
2944		3002
2945		3003
2946		3004
2947		3005
2948		3006
2949		3007
2950		3008
2951		3009
2952		3010
2953		3011
2954		3012
2955		3013
2956		3014
2957		3015
2958		3016