

318 A Related Works

319 Private inference has been a promising solution to protect both data and model privacy during deep
320 learning inference. In recent years, there has been an increasing amount of literature on efficient
321 private inference. According to the optimization technique, these works can be categorized into three
322 types, i.e., 1) protocol optimization; 2) network optimization; and 3) joint optimization.

323 In protocol optimization, ABY [7] provides a highly efficient conversion between arithmetic sharing,
324 boolean sharing and Yao’s sharing, and construct mixed protocols. As an extension, ABY3 [34]
325 switches back and forth between three secret sharing schemes using three-party computation (3PC).
326 CypTFlow2 [39] proposes a new protocol for secure and comparison and division which enables
327 efficient non-linear operations such as ReLU. SiRNN [38] further proposes 2PC protocols for
328 bitwidth extension, mixed-precision linear and non-linear operations. CrypTen [27] proposes a
329 software framework that provides a flexible machine learning focused API. More recently, SecFloat
330 [37] proposes the crypto-friendly precise functionalities to build a library for 32-bit single-precision
331 floating-point operations and math functions. These works lack consideration for neural network
332 architecture and has limited communication reduction.

333 In network optimization, DeepReDuce [24] proposes to manually remove ReLUs with a three-step
334 optimization pipeline. SNL [6] proposes ReLU-aware optimization that leverages gradient-based
335 NAS to selectively linearize a subset of ReLUs. CryptoNAS [18] uses ReLU budget as a proxy
336 and leverages NAS to tailor ReLUs. PolyMPCNet[36] and SAFENet [33] replace ReLUs with
337 MPC-friendly polynomial, while Sphynx [5] proposes an MPC-friendly ReLU-efficient micro-search
338 space. SENet [30] innovatively measures the ReLU importance via layer pruning sensitivity and
339 automatically optimize the network to meet the target ReLU budget. DeepReShape [24] finds that
340 wider networks are more ReLU-efficient than the deeper ones and designs ReLU-efficient baseline
341 networks with with FLOPs-ReLU-Accuracy balance. Network optimization mainly focuses on ReLU
342 reduction which dominates the online communication, but total communication including convolution
343 and truncation cannot be optimized.

344 Unluckily, only using either protocol or network optimization just leads to limited efficiency improve-
345 ment. Delphi [44] jointly optimizes cryptographic protocols and network by gradually replacing
346 ReLU with quadratic approximation. COINN [23] simultaneously optimizes quantized network and
347 protocols with ciphertext-aware quantization and automated bitwidth configuration. Recently, [16]
348 proposes to use Winograd convolution for reducing the number of multiplications and design the
349 efficient convolution operation to reduce the communication cost. However, it does not take private
350 inference into consideration for Winograd algorithm, and still suffers tremendous communication
351 overhead. In this work, we jointly optimize the network and protocol and fully consider their coupling
352 properties.

353 B Details of Experiment Setup

354 **Private inference framework** CoPriv adopts CypTFlow2 [39] protocol for private inference. We
355 leverage the Athos [39] tool chain to convert both input and weight into fixed-point with the bit-width
356 41 and scale 12. We measure the communication and latency under a LAN setting [39] with 377
357 MBps bandwidth and 80ms echo latency. All of our experiments are evaluated on the Intel Xeon
358 Gold 5220R CPU @ 2.20GHz.

359 **Implementation of Winograd-based convolution protocol** The convolution protocol with Wino-
360 grad transformation and optimization is implemented in C++ with Eigen and Armadillo matrix
361 calculation library [41] in the CypTFlow2 [39] framework. We implement $F(2 \times 2, 3 \times 3)$ and
362 $F(4 \times 4, 3 \times 3)$ transformation for convolution with stride of 1 and $F(2 \times 2, 3 \times 3)$ transformation
363 when stride is 2. For CIFAR-100 dataset, we use $F(2 \times 2, 3 \times 3)$ transformation as the image resolu-
364 tion is small and for ImageNet dataset, we use $F(4 \times 4, 3 \times 3)$. We only apply $F(2 \times 2, 3 \times 3)$ for
365 stride of 2 on ImageNet dataset. When evaluating CoPriv, we determine the optimal sender according
366 to the analysis in Table 3 before inference. Winograd implementation enables us to measure the
367 communication cost and latency of each convolution module.

368 **Networks and datasets** We apply our proposed CoPriv to the widely used lightweight mobile
369 network MobileNetV2 [42] with different width multipliers, e.g., 0.75, 1.0 and 1.4 to trade off the

370 model accuracy and efficiency. We evaluate the top-1 accuracy and online and total communication
 371 on both CIFAR-100 and ImageNet dataset.

372 **Differentiable pruning and finetuning setups** We first search and prune redundant ReLUs for 90
 373 epochs and then finetune the pruned network for 180 epochs with SGD optimizer, cosine learning
 374 scheduler and 0.1 initial learning rate. We train our proposed CoPriv with self-distillation.

375 C Network Re-Parameterization Algorithm

376 Network/Structural re-parameterization is a useful technique proposed by RepVGG [13], and is
 377 extended to [10, 9, 12, 15, 11]. The core idea of re-parameterization is to decouple the training-time
 378 architecture (with high performance and low efficiency) and inference-time network architecture
 379 (with high efficiency). Re-parameterization is realized by converting one architecture to another via
 380 equivalently merging parameters together. Therefore, during inference time, the network architecture
 381 is not only efficient but also has the same high performance as the training-time architecture.

382 In this work, we can also leverage this technique to merge adjacent convolutions together after
 383 ReLU removal. For the network re-parameterization mentioned in Section 4.2, here we provide the
 384 following detailed algorithm 1 to equivalently merge the inverted residual block into a single dense
 385 convolution as shown in Figure 3. With the help of network re-parameterization, we further optimize
 386 the total communication including convolution and truncation.

Algorithm 1: Network Re-parameterization for Inverted Residual Block

Input : An inverted residual block with weights $W_{1 \times 1}$, $W_{3 \times 3}$, and $W'_{1 \times 1}$. The number of input and
 output channels N_{in}, N_{out} . The size of re-parameterized weights r .
Output : Regular convolution with re-parameterized weights W_r .

```

1  $W_r = \text{torch.eye}(N_{in});$ 
2  $W_r = W_r.\text{unsqueeze}(2).\text{unsqueeze}(2);$ 
3  $W_r = \text{torch.nn.functional.pad}(W_r, \text{pad}=(\frac{r-1}{2}, \frac{r-1}{2}, \frac{r-1}{2}, \frac{r-1}{2}));$ 
4  $W_r = \text{torch.nn.functional.conv2d}(W_r, W_{1 \times 1});$ 
5  $W_r = \text{torch.nn.functional.conv2d}(W_r, W_{3 \times 3}, \text{padding}=\frac{r-1}{2});$ 
6  $W_r = \text{torch.nn.functional.conv2d}(W_r, W'_{1 \times 1});$ 
7  $W_{res} = \text{torch.zeros}(N_{out}, N_{in}, r, r);$ 
8 for  $i \in [0, \dots, N_{out} - 1]$  do
9    $W_{res}[i, i, \lfloor r/2 \rfloor, \lfloor r/2 \rfloor] = 1;$ 
10  $W_r = W_r + W_{res};$ 
11 return  $W_r;$ 

```

387 D Details of Winograd Convolution

388 D.1 Comparison between Regular Convolution and Winograd Convolution

389 To help readers better understand the multiplication reduction of Winograd convolution, we demon-
 390 strate regular convolution and Winograd convolution in Figure 10. Given an input $I \in \mathbb{R}^{4 \times 4}$ and
 391 a filter $F \in \mathbb{R}^{3 \times 3}$, regular convolution requires $9 \times 4 = 36$ times multiplications (implemented
 392 using GEMM with im2col algorithm [4]) while $F(2 \times 2, 3 \times 3)$ Winograd transformation only
 393 requires $16 \times 1 = 16$ times multiplications (EWMM), which achieves $2.25 \times$ reduction. Moreover,
 394 $F(4 \times 4, 3 \times 3)$ with a larger tile size, i.e., 6 can further achieve $4 \times$ multiplication reduction. The
 395 improvement gets benefit from the Winograd’s ability to convert im2col to EWMM and calculate the
 396 whole tile in Winograd domain at once.

397 D.2 Details of Input Tiling and Padding

398 Given a large 2D input $I \in \mathbb{R}^{l \times l}$, where $l > m + r - 1$, the core technique for ensuring the
 399 equivalence of regular convolution and Winograd convolution is input tiling and padding. The output
 400 size $l' = l - r + 1$, the input tile size $n = m + r - 1$ and the total tile number T per channel is

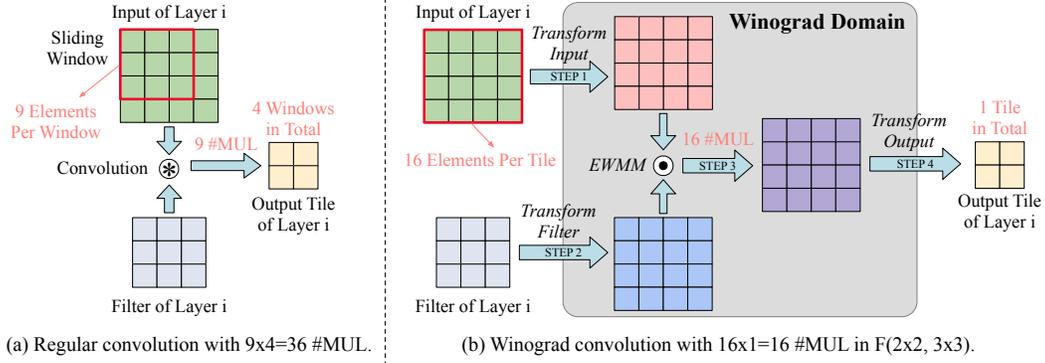


Figure 10: Comparison between (a) regular convolution and (b) Winograd convolution.

401 computed as

$$T = \lceil \frac{l'}{n} \rceil^2 = \lceil \frac{l-r+1}{m+r-1} \rceil^2,$$

402 where $\lceil \cdot \rceil$ denotes taking the upper bound value. For each tile, Winograd convolution is individually
 403 performed and results an output tile with $m \times m$ size. After all the tiles are computed with Winograd
 404 convolution, the output tiles are concatenated together to form the final output.

405 For some input size, the input cannot be covered by tiles. For instance, when leveraging $F(2 \times 2, 3 \times 3)$
 406 on the input $I \in \mathbb{R}^{7 \times 7}$, the rightmost and bottom pixels cannot be divided into a complete tile. To
 407 solve this problem, we pad these positions with 0 to enable the tiles totally cover the whole input.
 408 The correctness and equivalence can be proved with Eq. 1. Also, [16] shows the overhead caused by
 409 padding is negligible.

410 D.3 Support for Stride of 2 Winograd Convolution

411 Conventional Winograd convolution only supports stride $s = 1$ convolution filter. However, in recent
 412 efficient neural networks, e.g., MobileNetV2, EfficientNet has several stride of 2 layers to reduce
 413 the feature map size by half. To enable extreme optimization for efficient networks, we introduce
 414 $F(2 \times 2, 3 \times 3)$ for stride of 2 Winograd convolution for private inference.

415 There are various methods to construct stride of 2 Winograd kernel such as dividing input and
 416 convolution filter into different groups [46]. However, it is not a simple way to implement stride of 2
 417 Winograd kernel. [21] is an extremely convenient method using unified transformation matrices.

418 Based on [21], even positions of input and filter are computed by $F(2, 2)$ while odd positions
 419 are computed by regular convolution. Transformation matrices are derived as follows and can be
 420 computed using Eq. 1:

$$B^T = \begin{bmatrix} 1 & 0 & -1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & -1 & 0 & 1 \end{bmatrix}, \quad G = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad A^T = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \end{bmatrix}.$$

421 **Correctness analysis.** Here, we take a 1D algorithm as an example to prove the correctness
 422 Winograd convolution for stride of 2. The algorithm can be nested with itself to obtain a 2D algorithm
 423 [31].

424 Given input X and filter F as

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, \quad F = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \end{bmatrix}, \quad Y = X \otimes F = \begin{bmatrix} z_0 \\ z_1 \end{bmatrix}.$$

425 First, we calculate regular convolution with stride of 2 using im2col algorithm [4] as

$$Y_1 = \begin{bmatrix} x_0 & x_1 & x_2 \\ x_2 & x_3 & x_4 \end{bmatrix} \cdot \begin{bmatrix} y_0 \\ y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x_0y_0 + x_1y_1 + x_2y_2 \\ x_2y_0 + x_3y_1 + x_4y_2 \end{bmatrix},$$

426 thus, $z_0 = x_0y_0 + x_1y_1 + x_2y_2$ and $z_1 = x_2y_0 + x_3y_1 + x_4y_2$.

427 Then, we calculate Winograd convolution for stride of 2 as

$$Y = A^\top \cdot [(GF) \odot (B^\top X)],$$

428 and then

$$Y_2 = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \end{bmatrix} \cdot \left(\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} y_0 \\ y_1 \\ y_2 \end{bmatrix} \right) \odot \left(\begin{bmatrix} 1 & 0 & -1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & -1 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \right),$$

429 and further simplify the calculation as

$$Y_2 = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \end{bmatrix} \cdot \left(\begin{bmatrix} y_0 \\ y_1 \\ y_0 + y_2 \\ y_1 \\ y_2 \end{bmatrix} \right) \odot \left(\begin{bmatrix} x_0 - x_2 \\ x_1 \\ x_2 \\ x_3 \\ x_4 - x_2 \end{bmatrix} \right) = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} x_0y_0 - x_2y_0 \\ x_1y_1 \\ x_2y_0 + x_2y_2 \\ x_3y_1 \\ x_4y_2 - x_2y_2 \end{bmatrix},$$

430 therefore, the convolution result is

$$Y_2 = \begin{bmatrix} x_0y_0 + x_1y_1 + x_2y_2 \\ x_2y_0 + x_3y_1 + x_4y_2 \end{bmatrix} = Y_1.$$

431 D.4 Transformation Matrices for Winograd Convolution

432 We provide the transformation matrices A, B, G for $F(2 \times 2, 3 \times 3)$ and $F(4 \times 4, 3 \times 3)$ Winograd
433 transformation based on polynomial Chinese remainder theorem (CRT) or Lagrange interpolation
434 [31].

435 For $F(2 \times 2, 3 \times 3)$, we have

$$B^\top = \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix}, \quad G = \begin{bmatrix} 1 & 0 & 0 \\ 1/2 & 1/2 & 1/2 \\ 1/2 & -1/2 & 1/2 \\ 0 & 0 & 1 \end{bmatrix}, \quad A^\top = \begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 1 & -1 & -1 \end{bmatrix}.$$

436 For $F(4 \times 4, 3 \times 3)$, we have

$$B^\top = \begin{bmatrix} 4 & 0 & -5 & 0 & 1 & 0 \\ 0 & -4 & -4 & 1 & 1 & 0 \\ 0 & 4 & -4 & -1 & 1 & 0 \\ 0 & -2 & -1 & 2 & 1 & 0 \\ 0 & 2 & -1 & -2 & 1 & 0 \\ 0 & 4 & 0 & -5 & 0 & 1 \end{bmatrix}, \quad G = \begin{bmatrix} 1/4 & 0 & 0 \\ -1/6 & -1/6 & -1/6 \\ -1/6 & 1/6 & -1/6 \\ 1/24 & 1/12 & 1/6 \\ 1/24 & -1/12 & 1/6 \\ 0 & 0 & 1 \end{bmatrix},$$

$$A^\top = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 0 \\ 0 & 1 & -1 & 2 & -2 & 0 \\ 0 & 1 & 1 & 4 & 4 & 0 \\ 0 & 1 & -1 & 8 & -8 & 1 \end{bmatrix}.$$

438 The correctness analysis is the same with Section D.3.

References

- 439
- 440 [1] Syed Asad Alam, Andrew Anderson, Barbara Barabasz, and David Gregg. Winograd convo-
441 lution for deep neural networks: Efficient point selection. *ACM Transactions on Embedded*
442 *Computing Systems*, 21(6):1–28, 2022.
- 443 [2] Barbara Barabasz, Andrew Anderson, Kirk M Soodhalter, and David Gregg. Error analysis and
444 improving the accuracy of winograd convolution for deep neural networks. *ACM Transactions*
445 *on Mathematical Software (TOMS)*, 46(4):1–33, 2020.
- 446 [3] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients
447 through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- 448 [4] Kumar Chellapilla, Sidd Puri, and Patrice Simard. High performance convolutional neural
449 networks for document processing. In *Tenth international workshop on frontiers in handwriting*
450 *recognition*. Suvisoft, 2006.
- 451 [5] Minsu Cho, Zahra Ghodsi, Brandon Reagen, Siddharth Garg, and Chinmay Hegde. Sphynx: A
452 deep neural network design for private inference. *IEEE Security & Privacy*, 20(5):22–34, 2022.
- 453 [6] Minsu Cho, Ameya Joshi, Brandon Reagen, Siddharth Garg, and Chinmay Hegde. Selective
454 network linearization for efficient private inference. In *International Conference on Machine*
455 *Learning*, pages 3947–3961. PMLR, 2022.
- 456 [7] Daniel Demmler, Thomas Schneider, and Michael Zohner. Aby-a framework for efficient
457 mixed-protocol secure two-party computation. In *NDSS*, 2015.
- 458 [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
459 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*
460 *recognition*, pages 248–255. Ieee, 2009.
- 461 [9] Xiaohan Ding, Honghao Chen, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Repmlpnet:
462 Hierarchical vision mlp with re-parameterized locality. In *Proceedings of the IEEE/CVF*
463 *Conference on Computer Vision and Pattern Recognition*, pages 578–587, 2022.
- 464 [10] Xiaohan Ding, Yuchen Guo, Guiguang Ding, and Jungong Han. Acnet: Strengthening the
465 kernel skeletons for powerful cnn via asymmetric convolution blocks. In *Proceedings of the*
466 *IEEE/CVF international conference on computer vision*, pages 1911–1920, 2019.
- 467 [11] Xiaohan Ding, Tianxiang Hao, Jianchao Tan, Ji Liu, Jungong Han, Yuchen Guo, and Guiguang
468 Ding. Resrep: Lossless cnn pruning via decoupling remembering and forgetting. In *Proceedings*
469 *of the IEEE/CVF International Conference on Computer Vision*, pages 4510–4520, 2021.
- 470 [12] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Diverse branch block:
471 Building a convolution as an inception-like unit. In *Proceedings of the IEEE/CVF Conference*
472 *on Computer Vision and Pattern Recognition*, pages 10886–10895, 2021.
- 473 [13] Xiaohan Ding, Xiangyu Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun.
474 Repvgg: Making vgg-style convnets great again. In *Proceedings of the IEEE/CVF conference*
475 *on computer vision and pattern recognition*, pages 13733–13742, 2021.
- 476 [14] Javier Fernandez-Marques, Paul Whatmough, Andrew Mundy, and Matthew Mattina. Searching
477 for winograd-aware quantized networks. *Proceedings of Machine Learning and Systems*,
478 2:14–29, 2020.
- 479 [15] Yonggan Fu, Haichuan Yang, Jiayi Yuan, Meng Li, Cheng Wan, Raghuraman Krishnamoorthi,
480 Vikas Chandra, and Yingyan Lin. Depthshrinker: a new compression paradigm towards boosting
481 real-hardware efficiency of compact neural networks. In *International Conference on Machine*
482 *Learning*, pages 6849–6862. PMLR, 2022.
- 483 [16] Vinod Ganesan, Anwesh Bhattacharya, Pratyush Kumar, Divya Gupta, Rahul Sharma, and
484 Nishanth Chandran. Efficient ml models for practical secure inference. *arXiv preprint*
485 *arXiv:2209.00411*, 2022.

- 486 [17] Karthik Garimella, Zahra Ghodsi, Nandan Kumar Jha, Siddharth Garg, and Brandon Reagen.
487 Characterizing and optimizing end-to-end systems for private inference. In *ACM International*
488 *Conference on Architectural Support for Programming Languages and Operating Systems*
489 *(ASPLOS)*, ASPLOS 2023, page 89–104, New York, NY, USA, 2023. Association for Computing
490 Machinery.
- 491 [18] Zahra Ghodsi, Akshaj Kumar Veldanda, Brandon Reagen, and Siddharth Garg. Cryptonas:
492 Private inference on a relu budget. *Advances in Neural Information Processing Systems*,
493 33:16961–16971, 2020.
- 494 [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
495 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
496 pages 770–778, 2016.
- 497 [20] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan,
498 Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3.
499 In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1314–1324,
500 2019.
- 501 [21] Chengcheng Huang, Xiaoxiao Dong, Zhao Li, Tengting Song, Zhenguo Liu, and Lele Dong.
502 Efficient stride 2 winograd convolution method using unified transformation matrices on fpga.
503 In *2021 International Conference on Field-Programmable Technology (ICFPT)*, pages 1–9.
504 IEEE, 2021.
- 505 [22] Zhicong Huang, Wen-jie Lu, Cheng Hong, and Jiansheng Ding. Cheetah: Lean and fast secure
506 {Two-Party} deep neural network inference. In *31st USENIX Security Symposium (USENIX*
507 *Security 22)*, pages 809–826, 2022.
- 508 [23] Siam Umar Hussain, Mojan Javaheripi, Mohammad Samragh, and Farinaz Koushanfar. Coinn:
509 Crypto/ml codesign for oblivious inference via neural networks. In *Proceedings of the 2021*
510 *ACM SIGSAC Conference on Computer and Communications Security*, pages 3266–3281, 2021.
- 511 [24] Nandan Kumar Jha, Zahra Ghodsi, Siddharth Garg, and Brandon Reagen. Deepreduce: Relu
512 reduction for fast private inference. In *International Conference on Machine Learning*, pages
513 4839–4849. PMLR, 2021.
- 514 [25] Nandan Kumar Jha and Brandon Reagen. Deepreshape: Redesigning neural networks for
515 efficient private inference. *arXiv preprint arXiv:2304.10593*, 2023.
- 516 [26] Chiraag Juvekar, Vinod Vaikuntanathan, and Anantha Chandrakasan. GAZELLE: A low
517 latency framework for secure neural network inference. In *27th {USENIX} Security Symposium*
518 *({USENIX} Security 18)*, pages 1651–1669, 2018.
- 519 [27] Brian Knott, Shobha Venkataraman, Awni Hannun, Shubho Sengupta, Mark Ibrahim, and
520 Laurens van der Maaten. Crypten: Secure multi-party computation meets machine learning.
521 *Advances in Neural Information Processing Systems*, 34:4961–4973, 2021.
- 522 [28] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
523 2009.
- 524 [29] Nishant Kumar, Mayank Rathee, Nishanth Chandran, Divya Gupta, Aseem Rastogi, and Rahul
525 Sharma. Cryptflow: Secure tensorflow inference. In *IEEE Symposium on Security and Privacy*
526 *(SP)*, pages 336–353. IEEE, 2020.
- 527 [30] Souvik Kundu, Shunlin Lu, Yuke Zhang, Jacqueline Liu, and Peter A Beerel. Learning
528 to linearize deep neural networks for secure and efficient private inference. *arXiv preprint*
529 *arXiv:2301.09254*, 2023.
- 530 [31] Andrew Lavin and Scott Gray. Fast algorithms for convolutional neural networks. In *Proceedings*
531 *of the IEEE conference on computer vision and pattern recognition*, pages 4013–4021, 2016.
- 532 [32] Zechun Liu, Haoyuan Mu, Xiangyu Zhang, Zichao Guo, Xin Yang, Kwang-Ting Cheng, and
533 Jian Sun. Metapruning: Meta learning for automatic neural network channel pruning. In
534 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3296–3305,
535 2019.

- 536 [33] Qian Lou, Yilin Shen, Hongxia Jin, and Lei Jiang. Safenet: A secure, accurate and fast neural
537 network inference. In *International Conference on Learning Representations*, 2021.
- 538 [34] Payman Mohassel and Peter Rindal. Aby3: A mixed protocol framework for machine learning.
539 In *ACM SIGSAC conference on computer and communications security*, pages 35–52, 2018.
- 540 [35] Payman Mohassel and Yupeng Zhang. Secureml: A system for scalable privacy-preserving
541 machine learning. In *2017 IEEE symposium on security and privacy (SP)*, pages 19–38. IEEE,
542 2017.
- 543 [36] Hongwu Peng, Shanglin Zhou, Yukui Luo, Shijin Duan, Nuo Xu, Ran Ran, Shaoyi Huang,
544 Chenghong Wang, Tong Geng, Ang Li, et al. Polymcnet: Towards relu-free neural architecture
545 search in two-party computation based private inference. *arXiv preprint arXiv:2209.09424*,
546 2022.
- 547 [37] Deevashwer Rathee, Anwesh Bhattacharya, Rahul Sharma, Divya Gupta, Nishanth Chandran,
548 and Aseem Rastogi. Secfloat: Accurate floating-point meets secure 2-party computation. In
549 *2022 IEEE Symposium on Security and Privacy (SP)*, pages 576–595. IEEE, 2022.
- 550 [38] Deevashwer Rathee, Mayank Rathee, Rahul Kranti Kiran Goli, Divya Gupta, Rahul Sharma,
551 Nishanth Chandran, and Aseem Rastogi. Sirnn: A math library for secure rnn inference. In
552 *2021 IEEE Symposium on Security and Privacy (SP)*, pages 1003–1020. IEEE, 2021.
- 553 [39] Deevashwer Rathee, Mayank Rathee, Nishant Kumar, Nishanth Chandran, Divya Gupta, Aseem
554 Rastogi, and Rahul Sharma. Cryptflow2: Practical 2-party secure inference. In *Proceedings of*
555 *the 2020 ACM SIGSAC Conference on Computer and Communications Security*, pages 325–342,
556 2020.
- 557 [40] Brandon Reagen, Woo-Seok Choi, Yeongil Ko, Vincent T Lee, Hsien-Hsin S Lee, Gu-Yeon Wei,
558 and David Brooks. Cheetah: Optimizing and accelerating homomorphic encryption for private
559 inference. In *IEEE International Symposium on High-Performance Computer Architecture*
560 *(HPCA)*, pages 26–39. IEEE, 2021.
- 561 [41] Conrad Sanderson and Ryan Curtin. Armadillo: a template-based c++ library for linear algebra.
562 *Journal of Open Source Software*, 1(2):26, 2016.
- 563 [42] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen.
564 Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference*
565 *on computer vision and pattern recognition*, pages 4510–4520, 2018.
- 566 [43] Liyan Shen, Ye Dong, Binxing Fang, Jinqiao Shi, Xuebin Wang, Shengli Pan, and Ruisheng
567 Shi. Abnn2: secure two-party arbitrary-bitwidth quantized neural network predictions. In
568 *Proceedings of the 59th ACM/IEEE Design Automation Conference*, pages 361–366, 2022.
- 569 [44] Wenting Zheng Srinivasan, PMRL Akshayaram, and Popa Raluca Ada. Delphi: A cryptographic
570 inference service for neural networks. In *Proc. 29th USENIX Secur. Symp*, pages 2505–2522,
571 2019.
- 572 [45] Kevin Vincent, Kevin Stephano, Michael Frumkin, Boris Ginsburg, and Julien Demouth. On
573 improving the numerical stability of winograd convolutions. 2017.
- 574 [46] Juan Yopez and Seok-Bum Ko. Stride 2 1-d, 2-d, and 3-d winograd for convolutional neural
575 networks. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 28(4):853–863,
576 2020.