
A Generic Hybrid 2PC Framework with Application to Private Inference of Unmodified Neural Networks (Extended Abstract)

Lennart Braun
Aarhus University
braun@cs.au.dk

Rosario Cammarota
Intel Labs
rosario.cammarota@intel.com

Thomas Schneider
TU Darmstadt
schneider@encrypto.cs.tu-darmstadt.de

Abstract

1 We present a new framework for generic mixed-protocol secure two-party computa-
2 tion (2PC) and private evaluation of neural networks based on the recent MOTION
3 framework (Braun et al., ePrint '20). We implement five different 2PC protocols in
4 the semi-honest setting – Yao’s garbled circuits, arithmetic and Boolean variants
5 of Goldreich-Micali-Wigderson (GMW), and two secret-sharing-based protocols
6 from ABY2.0 (Patra et al., USENIX Security '21) – together with 20 conversions
7 among each other and new optimizations. We explore the feasibility of evaluating
8 neural networks with 2PC without making modifications to their structure, and
9 provide secure tensor data types and specialized building blocks for common tensor
10 operations. By supporting the Open Neural Network Exchange (ONNX) file format,
11 this yields an easy-to-use solution for privately evaluating neural networks, and is
12 interoperable with industry-standard deep learning frameworks such as TensorFlow
13 and PyTorch. By exploiting the networks’ high-level structure and using common
14 2PC techniques, we obtain a performance that is comparable to that of recent,
15 highly optimized works and significantly better than when using generic 2PC for
16 low-level hybrid circuits.

17 1 Introduction

18 **Secure Computation** Secure multi-party computation (MPC) allows a set of mutually distrusting
19 parties to evaluate a functionality on their private inputs and obtain private outputs such that each
20 party learns only what it can deduce from its own input and output, [21]. The two-party case is
21 also designated as secure two-party computation (2PC). MPC allows computation on data which
22 the owners are either reluctant or not allowed to share with each other, e.g., due to data privacy
23 regulations. Use cases include computing statistics on sensitive medical and corporate data [14], or
24 outsourcing of computation to cloud services [25]. The concept of secure computation dates back to
25 the 1980s [47, 46, 18], and has become increasingly practical over the last two decades [34, 39].

26 **Machine Learning** As a completely different technology, machine learning techniques such as
27 deep neural networks have been applied to many different tasks over the last years, e.g., image classi-
28 fication [19]. Since training neural networks requires large datasets and significant computational

29 resources, companies might be inclined to keep the resulting models secret to protect their investment.
30 On the other hand, customers may be hesitant to just pass their input to the service provider.

31 **Privacy-Preserving Machine Learning** The combination of multi-party computation and machine
32 learning allows a variety of privacy-preserving online services [32, 23, 4, 8]. For example, to classify
33 a picture, a customer can run an MPC protocol with a service provider instead of just uploading
34 the picture. It then gets the result of the classification without having to reveal the input, while the
35 provider can keep its model secret.

36 **Our Contributions** In this work, we present the first secure two-party computation framework
37 combining *five different protocols* such that shared data can be converted arbitrarily among the
38 protocols. In the framework, we combine Yao’s protocol [46, 31], the arithmetic and Boolean variants
39 of GMW [18] and the secret-sharing-based arithmetic and Boolean protocols from ABY2.0 [38]. We
40 analyze the protocols, and introduce new *optimizations* and conversion protocols. We integrate the
41 protocols into the recent MPC framework MOTION [7] which is partially redesigned and extended.
42 Furthermore, we explore the feasibility of private neural network inference using *standard MPC*
43 *techniques* without making modification to the networks. To this end, we discuss protocols for
44 common neural network operations, and implement them as specialized building blocks in MOTION.
45 By supporting the *ONNX* file format [37] we enable *interoperability with deep learning frameworks*
46 *used in industry* such as TensorFlow and PyTorch. Finally, we evaluate the performance of our
47 framework using various benchmarks, and compare the implemented protocols among themselves
48 and with prior work. The specialized building blocks perform clearly better for neural networks than
49 the equivalent generic 2PC protocols for hybrid circuits, and we achieve a performance comparable
50 to recent, highly optimized works such as GAZELLE [24] and DELPHI [35].

51 **Related Work** The most relevant works regarding generic MPC are ABY [13], MOTION [7], and
52 ABY2.0 [38] upon which we heavily build (see § 2, § 3). TASTY [22], ABY and MOTION are
53 software frameworks for 2PC and MPC in the semi-honest security model supporting Boolean and
54 arithmetic operations. MP-SPDZ [26] and SCALE-MAMBA [2] provide protocols that are secure
55 in the malicious model. EzPC [10] and HyCC [9] can compile high-level function descriptions into
56 hybrid circuits usable for MPC. There is a large body of work on private evaluation of neural networks,
57 e.g. [4, 16, 36, 24, 42, 41]. Many of these modify the networks, e.g., by quantizing weights, to obtain
58 more efficient protocols, whereas we try to preserve the original network as much as possible.

59 **Organization** In § 2, we give an overview of the used protocols. § 3 covers our improvements
60 to the MOTION framework, and in § 4, we discuss the results of our experimental evaluation.
61 Supplementary material can be found in the appendix.

62 2 Protocols

63 In this work, we consider *five* generic 2PC protocols for Boolean circuits (denoted with Y, B, β) and
64 arithmetic circuits over rings \mathbb{Z}_{2^e} (denoted with A, α). The protocols are secure in the semi-honest
65 security setting, i.e., corrupted parties follow the protocol, but try to learn additional information. As
66 in ABY [13], we use Yao’s garbled circuit protocol (Y) [47, 46] and the arithmetic (A) and Boolean
67 (B) variants of the GMW protocol [18]. Additionally, we combine these with the new arithmetic
68 (α) and Boolean (β) protocols introduced in ABY2.0 [38]. See § A.2 for a short summary of the
69 protocols. For all protocol operations, we consider single instruction multiple data (SIMD) variants.
70 They operate element-wise on vectors of values, and can be implemented much more efficiently.

71 **Conversions** Some operations (e.g. additions, multiplications) are naturally expressed as arithmetic
72 circuits, and others are more efficiently represented in Boolean circuits (e.g. comparisons). Hence, it
73 is often advantageous to combine both kinds into a hybrid circuit. To evaluate such a hybrid circuit
74 with 2PC protocols, we need to convert between different representations. We provide *conversions*
75 *between all five protocols*. Many are based on prior work [13, 38], some are new or optimized.

76 **Neural Networks** We use the generic 2PC protocols to securely evaluate the *tensor operations*
77 *in neural networks*. Here, we encode the values in \mathbb{Z}_{2^t} using a fixed-point representation and the
78 truncation protocol by [36] instead of using floating-point numbers. Since some of the operations are
79 more efficiently computed with an arithmetic protocol and others with a Boolean protocol, we use the
80 conversions to change the data representation as necessary. For fully-connected and convolutional
81 layers, we use generalizations of the integer multiplication in the A and α protocols (cf. [36]).
82 AvgPool is a linear operation and needs only a truncation. MaxPool is based on optimized Boolean
83 circuits. For the ReLU activation function we provide different variants, e.g., based on special
84 bit-integer multiplication protocols. More details about these tensor operations are given in § A.3.

85 **ABY vs. ABY2.0** We compare the ABY [13] (A, B, Y) and ABY2.0 [38] (α, β, Y) protocol suites:
86 Storing an α shares requires twice the space compared to an A share. Linear operations are computed
87 in both cases locally without interaction, although ABY2.0 requires more (local) computation during
88 the setup phase. The ABY2.0 *multiplication* needs only half of the online communication compared
89 to GMW [38]. More generally, the online communication in GMW is linear in the size of the inputs,
90 whereas for ABY2.0, it is linear in the output size. This also holds when the multiplication protocols
91 are generalized, e.g., to matrix multiplications or convolutions.

92 The *setup phase*, i.e., the part of the protocol that can be executed before the inputs are known,
93 of GMW depends only on the number of multiplications. For ABY2.0, it also depends on the
94 circuit structure, which makes it more costly in general. In the case of neural networks, the (matrix)
95 multiplication operations are relatively few and relatively large compared to integer multiplications
96 in a normal arithmetic circuit. Since MaxPool and ReLU layers repeatedly apply the same basic
97 function, the setup can be computed in batches with SIMD operations. Thus, the disadvantages of a
98 function-dependent setup phase do not carry much weight in the case of neural networks.

99 We compare the *conversions* among the A, B , and Y protocols to those among α, β , and Y . The
100 latter were presented by [38] who also showed that they compare favorably to the original ABY
101 conversions [13]. The newer conversions presented in our work improve on ABY [13] and the
102 differences have become smaller. Looking at costs, a pattern becomes apparent: The conversions
103 from A/B to Y require two rounds while the other direction is free of communication in the online
104 phase. On the other hand, converting from α/β to Y and vice versa costs one round. Overall, in a
105 deep circuit (e.g., a neural network) where different layers are computed alternately with Y and one
106 of A/α , the number of rounds in the online phase are the same.

107 We examine different ways to compute *ReLU layers* in neural networks. The online communication
108 costs for the ABY2.0 protocols are significantly lower compared to their GMW counterparts. The
109 total communication costs are either similar or slightly in favor of ABY2.0.

110 3 Extending MOTION for Neural Networks

111 **Extending the Framework** In this work, we build upon and extend the MOTION framework
112 for mixed-protocol multi-party computation [7]. It is not only a library implementing some MPC
113 protocols, but also a *framework* that consists of useful components which can be composed and
114 extended as needed. We implement the five generic 2PC protocols (§ 2) and conversions among them,
115 as well as building blocks for neural network evaluation.

116 We redesigned many of the framework’s interfaces to reduce overhead and *improve flexibility* by
117 decoupling the components (see § B for more details). Based on the terminology of circuits,
118 MOTION [7] implements the primitive operations of protocols in Gate classes, which it evaluates
119 using fibers, i.e., threads that run in user space. We extracted the execution code into Executor
120 classes that take a collection of gates and execute them according to some strategy. This improves the
121 framework’s flexibility by allowing the implementation of different strategies. As in [7], the default
122 executor for general-purpose MPC with arbitrary circuits creates fibers for all gates which are then
123 executed by a thread pool.

124 **MOTION for Neural Networks** For evaluating *neural networks*, we take a different approach:
125 Many neural networks have a simple or even straight-line topology and are composed from a small
126 set of different tensor operations (e.g., linear and pooling layers and activation functions). Hence,
127 instead of compiling their descriptions into low-level circuits, we provide optimized building blocks
128 that directly implement common high-level operations on shared tensors. The tensor operations
129 are more computationally complex than the primitive operations of a circuit, but they have a very
130 regular structure, e.g., a *single* activation function is applied to *all* elements of a tensor. We exploit
131 this knowledge and use a specialized executor to evaluate the tensor operations sequentially, while
132 parallelizing the operations themselves with multi-threading and SIMD operations.

133 We enable *interoperability* with industry standard deep learning frameworks such as TensorFlow and
134 PyTorch by supporting the ONNX file format [37], an open standard for deep learning models. Our
135 new `OnnxAdapter` based on the ONNX library¹ can parse the description of neural networks and
136 automatically constructs the respective tensor operations in MOTION. This makes it easy to privately
137 evaluate models built with common deep learning frameworks. Moreover, we built tools to display
138 information about the neural network and estimate the costs of a secure evaluation.

139 4 Performance Evaluation

140 We extensively benchmarked our framework and compare the performance with prior and concurrent
141 work. See § C.1 for more details in the context of generic 2PC.

142 **Neural Networks** We use the CryptoNets [16] and MiniONN [33] networks for the MNIST [30]
143 dataset to compare the performance of our neural network building blocks with generic 2PC of hybrid
144 circuits generated by HyCC [9]. In almost every case, our dedicated building blocks outperform the
145 generic 2PC implementations of our framework, ABY [13], and MOTION [7]. Detailed benchmark
146 results are given in § C.2.

147 To examine the performance on larger neural networks, we use the MiniONN [33] neural network
148 for the CIFAR-10 [29] dataset, and compare the online run-times using our neural network building
149 blocks with results for GAZELLE [24], DELPHI [35], and CryptFlow2 [41].² The results (see
150 § C.2) show that the recently published CryptFlow2 offers clearly the best performance, but the
151 performance with our building blocks is comparable to that of GAZELLE and DELPHI.

152 **ABY vs. ABY2.0** According to our experiments, SIMD seems to be *even more important* for a
153 good performance with the ABY2.0 secret-sharing-based protocols than for GMW. Without SIMD,
154 the function-dependent setup phase is a clear drawback of the ABY2.0 sharings, and can result in
155 significantly slower setup phases compared to GMW. With SIMD, the differences in the setup phase
156 diminish. Sometimes, the measured setup run-time was even lower for ABY2.0 compared to the
157 ABY protocols.

158 For the evaluation of circuits with the generic 2PC implementation, there is no clear winner in the
159 online phase: In most of our experiments, the ABY protocols are better in some settings, and the
160 ABY2.0 protocols are better in other settings. An exception is the CryptoNets [16] benchmark, where
161 the ABY2.0 protocols perform almost always worse than the ABY protocols in the online phase.

162 When evaluating a neural network using our specialized building blocks, the ABY2.0 protocols are
163 the best choice regarding the online run-times. This is observed best when benchmarking ReLU
164 operations and the two MiniONN [33] networks.

¹ONNX: <https://github.com/onnx/onnx>

²The latter results are taken from [24] and [41], and have been obtained in a different experimental setup and with different bit sizes.

165 **Acknowledgments**

166 This work is supported by the European Research Council (ERC) under the European Unions’s
167 Horizon 2020 research and innovation programme under grant agreements No. 803096 (SPEC) and
168 No. 850990 (PSOTI), and the Carlsberg Foundation under the Semper Ardens Research Project
169 CF18-112 (BCM). It was co-funded by the Deutsche Forschungsgemeinschaft (DFG) — SFB 1119
170 CROSSING/236615297 and GRK 2050 Privacy & Trust/251805230, and by the German Federal
171 Ministry of Education and Research and the Hessen State Ministry for Higher Education, Research
172 and the Arts within ATHENE.

173 **References**

- 174 [1] V. A. Abril, P. Maene, N. Mertens, D. Sijacic, and N. Smart. ‘Bristol Fashion’ MPC Circuits. 2019. URL:
175 <https://homes.esat.kuleuven.be/~nsmart/MPC/>.
- 176 [2] A. Aly, K. Cong, K. Koch, M. Keller, D. Rotaru, O. Scherer, P. Scholl, N. P. Smart, T. Tanguy, and
177 T. Wood. *SCALE-MAMBA Software*. <https://homes.esat.kuleuven.be/~nsmart/SCALE/>. 2018.
- 178 [3] G. Asharov, Y. Lindell, T. Schneider, and M. Zohner. *More Efficient Oblivious Transfer and Extensions
179 for Faster Secure Computation*. In: *Proceedings of the 2013 ACM SIGSAC Conference on Computer and
180 Communications Security, CCS 2013* (Berlin, Germany, Nov. 4–8, 2013). ACM, 2013.
- 181 [4] M. Barni, C. Orlandi, and A. Piva. *A privacy-preserving protocol for neural-network-based computation*.
182 In: *Proceedings of the 8th workshop on Multimedia & Security, MM&Sec 2006* (Geneva, Switzerland,
183 Sept. 26–27, 2006). ACM, 2006.
- 184 [5] D. Beaver. *Efficient Multiparty Protocols Using Circuit Randomization*. In: *Advances in Cryptology -
185 CRYPTO ’91, 11th Annual International Cryptology Conference, Proceedings* (Santa Barbara, CA, USA,
186 Aug. 11–15, 1991). Springer, 1992.
- 187 [6] M. Bellare, V. T. Hoang, S. Keelveedhi, and P. Rogaway. *Efficient Garbling from a Fixed-Key Blockcipher*.
188 In: *2013 IEEE Symposium on Security and Privacy, S&P 2013* (Berkeley, CA, USA, May 19–22, 2013).
189 IEEE, 2013.
- 190 [7] L. Braun, D. Demmler, T. Schneider, and O. Tkachenko. *MOTION – A Framework for Mixed-Protocol
191 Multi-Party Computation*. 2020. IACR Cryptology ePrint Archive: 2020/1137.
- 192 [8] J. Brickell, D. E. Porter, V. Shmatikov, and E. Witchel. *Privacy-preserving remote diagnostics*. In:
193 *Proceedings of the 2007 ACM Conference on Computer and Communications Security, CCS 2007*
194 (Alexandria, VA, USA, Oct. 28–31, 2007). ACM, 2007.
- 195 [9] N. Büscher, D. Demmler, S. Katzenbeisser, D. Kretzmer, and T. Schneider. *HyCC: Compilation of Hybrid
196 Protocols for Practical Secure Computation*. In: *Proceedings of the 2018 ACM SIGSAC Conference on
197 Computer and Communications Security, CCS 2018* (Toronto, ON, Canada, Oct. 15–19, 2018). ACM,
198 2018.
- 199 [10] N. Chandran, D. Gupta, A. Rastogi, R. Sharma, and S. Tripathi. *EzPC: Programmable and Efficient
200 Secure Two-Party Computation for Machine Learning*. In: *IEEE European Symposium on Security and
201 Privacy, EuroS&P 2019* (Stockholm, Sweden, June 17–19, 2019). IEEE, 2019.
- 202 [11] I. Damgård, M. Fitzi, E. Kiltz, J. B. Nielsen, and T. Toft. *Unconditionally Secure Constant-Rounds
203 Multi-party Computation for Equality, Comparison, Bits and Exponentiation*. In: *Theory of Cryptography,
204 Third Theory of Cryptography Conference, TCC 2006, Proceedings* (New York, NY, USA, Mar. 4–7,
205 2006). Springer, 2006.
- 206 [12] D. Demmler, G. Dessouky, F. Koushanfar, A. Sadeghi, T. Schneider, and S. Zeitouni. *Automated Synthesis
207 of Optimized Circuits for Secure Computation*. In: *Proceedings of the 22nd ACM SIGSAC Conference on
208 Computer and Communications Security, CCS 2015* (Denver, CO, USA, Oct. 12–16, 2015). ACM, 2015.
- 209 [13] D. Demmler, T. Schneider, and M. Zohner. *ABY - A Framework for Efficient Mixed-Protocol Secure
210 Two-Party Computation*. In: *22nd Annual Network and Distributed System Security Symposium, NDSS
211 2015* (San diego, CA, USA, Feb. 8–11, 2015). The Internet Society, 2015.
- 212 [14] W. Du and M. J. Atallah. *Secure Multi-Party Computation Problems and Their Applications: A Review
213 and Open Problems*. In: *Proceedings of the New Security Paradigms Workshop 2001* (Cloudcroft, NM,
214 USA, Sept. 10–13, 2001). ACM, 2001.
- 215 [15] S. Even, O. Goldreich, and A. Lempel. *A Randomized Protocol for Signing Contracts*. In: *Advances in
216 Cryptology: Proceedings of CRYPTO ’82* (Santa Barbara, CA, USA, Aug. 23–25, 1982). Plenum Press,
217 New York, 1983.
- 218 [16] R. Gilad-Bachrach, N. Dowlin, K. Laine, K. E. Lauter, M. Naehrig, and J. Wernsing. *CryptoNets:
219 Applying Neural Networks to Encrypted Data with High Throughput and Accuracy*. In: *Proceedings of the
220 33rd International Conference on Machine Learning, ICML 2016* (New York City, NY, USA, June 19–24,
221 2016). Vol. 48. JMLR Workshop and Conference Proceedings. JMLR.org, 2016.

- 222 [17] N. Gilboa. *Two Party RSA Key Generation*. In: *Advances in Cryptology - CRYPTO '99, 19th Annual International Cryptology Conference, Proceedings* (Santa Barbara, CA, USA, Aug. 15–19, 1999). Springer, 223 1999.
- 225 [18] O. Goldreich, S. Micali, and A. 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, STOC 1987* (New York, NY, USA, May 25–27, 1987). ACM, 1987.
- 226 [19] I. J. Goodfellow, Y. Bengio, and A. C. Courville. *Deep Learning*. MIT Press, 2016. URL: <http://www.deeplearningbook.org/>.
- 228 [20] C. Guo, J. Katz, X. Wang, and Y. Yu. *Efficient and Secure Multiparty Computation from Fixed-Key Block Ciphers*. In: *2020 IEEE Symposium on Security and Privacy, S&P 2020* (San Francisco, CA, USA, May 18, 2020–May 21, 2019). IEEE, 2020.
- 230 [21] C. Hazay and Y. Lindell. *Efficient Secure Two-Party Protocols. Techniques and Constructions*. Springer, 231 2010.
- 232 [22] W. Henecka, S. Kögl, A. Sadeghi, T. Schneider, and I. Wehrenberg. *TASTY: tool for automating secure two-party computations*. In: *Proceedings of the 17th ACM Conference on Computer and Communications Security, CCS 2010* (Chicago, IL, USA, Oct. 4–8, 2010). ACM, 2010.
- 233 [23] S. Jha, L. Kruger, and P. D. McDaniel. *Privacy Preserving Clustering*. In: *Computer Security - ESORICS 2005, 10th European Symposium on Research in Computer Security, Proceedings* (Milan, Italy, Sept. 12–14, 2005). Springer, 2005.
- 234 [24] C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan. *GAZELLE: A Low Latency Framework for Secure Neural Network Inference*. In: *27th USENIX Security Symposium, USENIX Security 2018* (Baltimore, MD, USA, Aug. 15–17, 2018). USENIX Association, 2018.
- 235 [25] S. Kamara, P. Mohassel, and M. Raykova. *Outsourcing Multi-Party Computation*. 2011. IACR Cryptology ePrint Archive: 2011/272.
- 236 [26] M. Keller. *MP-SPDZ: A Versatile Framework for Multi-Party Computation*. In: *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security, CCS 2020* (virtual conference, Nov. 9, 2020–Nov. 13, 2019). ACM, 2020.
- 237 [27] V. Kolesnikov and T. Schneider. *Improved Garbled Circuit: Free XOR Gates and Applications*. In: *Automata, Languages and Programming: 35th International Colloquium, ICALP 2008. Proceedings, Part II* (Reykjavik, Iceland, July 7–11, 2008). Springer, 2008.
- 238 [28] O. Kowalke. *Boost.Fiber*. Version 1.74.0. Library Documentation. 2020. URL: https://www.boost.org/doc/libs/1_74_0/libs/fiber/.
- 239 [29] A. Krizhevsky. *Learning Multiple Layers of Features from Tiny Images*. Tech. rep. <https://www.cs.utoronto.ca/~kriz/learning-features-2009-TR.pdf>. 2009.
- 240 [30] Y. LeCun, C. Cortes, and C. J. C. Burges. *The MNIST Database of Handwritten Digits*. 1998. URL: <http://yann.lecun.com/exdb/mnist/>.
- 241 [31] Y. Lindell and B. Pinkas. *A Proof of Security of Yao's Protocol for Two-Party Computation*. In: *Journal of Cryptology* 22.2 (2009).
- 242 [32] Y. Lindell and B. Pinkas. *Privacy Preserving Data Mining*. In: *Advances in Cryptology - CRYPTO 2000, 20th Annual International Cryptology Conference, Proceedings* (Santa Barbara, CA, USA, Aug. 20–24, 2000). Springer, 2000.
- 243 [33] J. Liu, M. Juuti, Y. Lu, and N. Asokan. *Oblivious Neural Network Predictions via MiniONN Transformations*. In: *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS 2017* (Dallas, TX, USA, Oct. 30–Nov. 3, 2017). ACM, 2017.
- 244 [34] D. Malkhi, N. Nisan, B. Pinkas, and Y. Sella. *Fairplay - Secure Two-Party Computation System*. In: *Proceedings of the 13th USENIX Security Symposium* (San Diego, CA, USA, Aug. 9–13, 2004). USENIX, 2004.
- 245 [35] P. Mishra, R. Lehmkuhl, A. Srinivasan, W. Zheng, and R. A. Popa. *Delphi: A Cryptographic Inference Service for Neural Networks*. In: *29th USENIX Security Symposium, USENIX Security 2020* (Boston, MA, USA, Aug. 12–14, 2020). USENIX Association, 2020.
- 246 [36] P. Mohassel and Y. Zhang. *SecureML: A System for Scalable Privacy-Preserving Machine Learning*. In: *2017 IEEE Symposium on Security and Privacy, S&P* (San Jose, CA, USA, May 22–26, 2017). IEEE, 2017.
- 247 [37] ONNX Project Contributors. *Open Neural Network Exchange (ONNX)*. 2019. URL: <https://onnx.ai>.
- 248 [38] A. Patra, T. Schneider, A. Suresh, and H. Yalame. *ABY2.0: Improved Mixed-Protocol Secure Two-Party Computation*. In: *30th USENIX Security Symposium, USENIX Security 2021* (Virtual Event, Aug. 13–15, 2021). USENIX Association, 2021.
- 249 [39] B. Pinkas, T. Schneider, N. P. Smart, and S. C. Williams. *Secure Two-Party Computation Is Practical*. In: *Advances in Cryptology - ASIACRYPT 2009, 15th International Conference on the Theory and Application of Cryptology and Information Security* (Tokyo, Japan, Dec. 6–10, 2009). Springer, 2009.

- 282 [41] D. Rathee, M. Rathee, N. Kumar, N. Chandran, D. Gupta, A. Rastogi, and R. Sharma. *CryptFlow2: Practical 2-Party Secure Inference*. In: *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security, CCS 2020* (virtual conference, Nov. 9, 2020–Nov. 13, 2019). ACM, 2020.
- 283
- 284
- 285 [42] M. S. Riazi, M. Samragh, H. Chen, K. Laine, K. E. Lauter, and F. Koushanfar. *XONN: XNOR-based Oblivious Deep Neural Network Inference*. In: *28th USENIX Security Symposium, USENIX Security 2019* (Santa Clara, CA, USA, Aug. 14–16, 2019). USENIX Association, 2019.
- 286
- 287
- 288 [43] M. Rosulek and L. Roy. *Three Halves Make a Whole? Beating the Half-Gates Lower Bound for Garbled Circuits*. In: *Advances in Cryptology - CRYPTO 2021 - 41st Annual International Cryptology Conference, CRYPTO 2021, Virtual Event, August 16-20, 2021, Proceedings, Part I*. Springer, 2021, pp. 94–124.
- 289
- 290
- 291 [44] S. Tillich and N. Smart. *(Bristol Format) Circuits of Basic Functions Suitable For MPC and FHE*. 2014. URL: <https://homes.esat.kuleuven.be/~nsmart/MPC/old-circuits.html>.
- 292
- 293 [45] X. Wang, A. J. Malozemoff, and J. Katz. *EMP-toolkit: Efficient MultiParty computation toolkit*. <https://github.com/emp-toolkit>. 2016.
- 294
- 295 [46] A. C.-C. Yao. *How to Generate and Exchange Secrets*. In: *27th Annual Symposium on Foundations of Computer Science, FOCS 1986* (Toronto, Canada, Oct. 27–29, 1986). IEEE Computer Society, 1986.
- 296
- 297 [47] A. C.-C. Yao. *Protocols for Secure Computation*. In: *23th Annual Symposium on Foundations of Computer Science, FOCS 1982* (Chicago, IL, USA, Nov. 3–5, 1982). IEEE Computer Society, 1982.
- 298
- 299 [48] S. Zahur, M. Rosulek, and D. Evans. *Two Halves Make a Whole - Reducing Data Transfer in Garbled Circuits Using Half Gates*. In: *Advances in Cryptology - EUROCRYPT 2015 - 34th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Proceedings, Part II* (Sofia, Bulgaria, Apr. 26–30, 2015). Springer, 2015.
- 300
- 301
- 302

303 Appendix A Protocols

304 Here, we give a high-level overview of the protocols covered and implemented in this work. The
 305 protocols are secure in the semi-honest security setting, i.e., corrupted parties follow the protocol,
 306 but try to learn additional information. We split the protocols into a setup and an online phase. The
 307 former is independent of the parties’ inputs and can be precomputed before the actual computation
 308 starts.

309 A.1 Auxiliary Protocols

310 **Secret Sharing** Secret sharing denotes a technique to split a secret into multiple shares such that it
 311 can be reconstructed from certain subsets of the shares. Important in our case is additive secret sharing
 312 over \mathbb{Z}_{2^ℓ} : To share $x \in \mathbb{Z}_{2^\ell}$, a dealer samples $x_1, \dots, x_N \in_R \mathbb{Z}_{2^\ell}$ uniformly at random such that
 313 $x = x_1 + \dots + x_N$, and distributes the x_i among the parties. This is denoted as $\langle x \rangle^A = (x_1, \dots, x_N)$.
 314 Given all shares the original secret can be recovered.

315 **Oblivious Transfer** Oblivious transfer (OT) [15] is a two-party protocol between a sender and a
 316 receiver. The sender inputs two messages m_0, m_1 and the receiver inputs a bit b . Then the receiver
 317 obtains m_b without learning m_{1-b} and the sender does not learn anything about b . There exist many
 318 variants of OT, e.g., correlated OT (C-OT) where the messages are chosen randomly under some
 319 correlation supplied by the sender [3].

320 **Multiplication Protocols** C-OT can be seen as computing a secret sharing of the product of two
 321 private bits. This has been extended to multiplication of integers, and generalized to products of
 322 vectors and matrices [17, 13, 36]. Moreover, two multiplication with private inputs can be combined
 323 to obtain a multiplication with two *shared* inputs. The arithmetic protocols covered below make
 324 heavy use of such OT-based multiplications.

325 A.2 General-Purpose 2PC

326 General-purpose 2PC protocols allow us to securely evaluate functionalities encoded in the form of
 327 Boolean or arithmetic circuits. Here we use and combine five existing protocols (denoted with $Y, A,$
 328 $B, \alpha,$ and β) with different properties.

329 **Yao’s Protocol (Y)** This is a 2PC protocol for Boolean circuits [47, 46]. One party, the *garbler*,
 330 creates an encoding of the circuit, the *garbled circuit*, and sends it to the second party, the *evaluator*.
 331 Given an encoding of the circuit inputs, the latter is able to obviously evaluate circuit to obtain an
 332 encoded output. This can be decoded with help of decoding information generated by the garbler.

333 To this end, the garbler creates for each circuit wire w two keys k_w^0, k_w^1 . For each gate g with
 334 input wires a, b and output wire c it creates a garbled table where for each possible combination of
 335 input values $\alpha, \beta \in \{0, 1\}$ the keys k_a^α and k_b^β are used to encrypt $k_c^{g(\alpha, \beta)}$. During evaluation the
 336 evaluator obtains k_w^b if the wire holds value b (denoted as $\langle b \rangle^Y$). Hence, given one key for each
 337 circuit inputs, the evaluator can evaluate the whole circuit while only seeing random-looking keys
 338 instead of the plain values. We implement Yao’s protocol with then state-of-the-art optimizations
 339 such as FreeXOR [27], fixed-key AES [6, 20], and Half-gates [48].³

340 **GMW (A/B)** The GMW protocol [18] is a generic MPC protocol for Boolean (B) and arithmetic
 341 (A) circuits in the full-threshold, semi-honest security setting. We consider GMW over the ring \mathbb{Z}_{2^ℓ} ,
 342 which for $\ell = 1$ is equivalent to bits $\{0, 1\}$ with \oplus and \wedge . The following description uses arithmetic
 343 notation, but unless stated otherwise holds also for the Boolean case.

344 A value $x \in \mathbb{Z}_{2^\ell}$ is shared using additive secret sharing $\langle x \rangle^A$. In the Boolean domain this is denoted
 345 as $\langle x \rangle^B$. Since the secret sharing scheme is linearly homomorphic, we can compute the sum of
 346 shared values by adding the shares locally. A multiplication $\langle z \rangle^A \leftarrow \langle x \rangle^A \cdot \langle y \rangle^A$ on the other hand
 347 requires interaction among the parties. As outlined in § A.1, shared values can be multiplied e.g. with
 348 OT-based multiplication protocols. Here, we use these to precompute multiplication triples (MTs) [5],
 349 i.e., random shared triples $(\langle a \rangle^A, \langle b \rangle^A, \langle c \rangle^A)$ such that $ab = c$, in the setup phase. Then, during the
 350 online phase the inputs x, y are masked with a, b from the MT, and then reconstructed. Finally a
 351 sharing of z can be computed using a linear combination of shared values with public coefficients.
 352 We also support mixed products of a bit $\langle b \rangle^B$ and a number $\langle n \rangle^A$.

353 **ABY2.0 Sharing (α/β)** Recently, ABY2.0 [38] was published as a successor protocol suite to the
 354 original ABY protocols [13]. They also combine Yao’s protocol (Y) with a Boolean and an arithmetic
 355 secret-sharing-based protocol. Instead of using GMW for the latter, they designed a new protocols
 356 which “uses a different perspective of [Beaver’s circuit randomization] technique” [38]. Furthermore,
 357 they designed a new set of conversion protocols and building blocks for various applications. To
 358 distinguish the new protocols, we use the term *ABY2.0 sharing* in this work and denote them with α
 359 and β .

360 To share a value $x \in \mathbb{Z}_{2^\ell}$, it is masked with a random value $\Delta_x \leftarrow x + \delta_x$ and the mask is additively
 361 secret shared among the parties: $\langle x \rangle^\alpha = (\Delta_x; \langle \delta_x \rangle^A)$. Note that δ_x can be generated independently
 362 of x during the setup phase. Addition works locally by adding the shares. For multiplications
 363 $\langle z \rangle^\alpha \leftarrow \langle x \rangle^\alpha \cdot \langle y \rangle^\alpha$, we need to compute $\langle \delta_x \cdot \delta_y \rangle^A \leftarrow \langle \delta_x \rangle^A \cdot \langle \delta_y \rangle^A$ in the setup phase. Then, in
 364 the online phase we can locally compute $\langle z \rangle^A$, and convert this to $\langle z \rangle^\alpha$ by masking it with δ_z and
 365 reconstructing the $\Delta_z = z + \delta_z$. Compared to GMW, only one reconstruction is needed, which
 366 halves the communication during the online phase.

367 We also support special products between shared bits and numbers $\langle b \cdot n \rangle^\alpha \leftarrow \langle b \rangle^\beta \cdot \langle n \rangle^\alpha$ and
 368 $\langle b_1 \cdot b_2 \rangle^\alpha \leftarrow \langle b_1 \rangle^\beta \cdot \langle b_2 \rangle^\beta$ and *improve* the setup phase compared to [38].

369 **Comparison** The three protocols have different properties: Yao’s protocol has a constant round
 370 complexity, whereas the other protocols need one round of interaction for each non-linear layer of the
 371 circuit. On the other hand transferring the garbled circuit needs a larger amount of communication.
 372 Hence, which approach performs better depends on the available network bandwidth and latency. For
 373 a more detailed comparison of the ABY and ABY2.0 protocol suites, see § 2.

374 **Conversion Protocols** To combine the different protocols and exploit their respective advantages,
 375 we provide conversion protocols between the five kinds of sharings. Figure 1 gives an overview of

³This was considered state-of-the-art until the recently published [43].

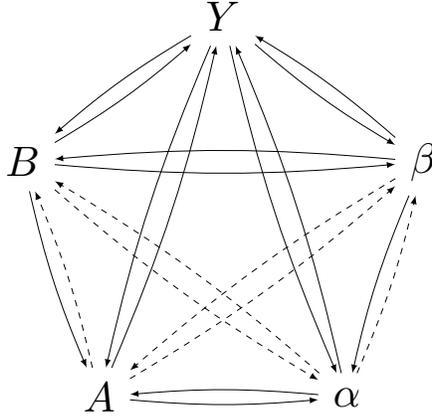


Figure 1: Overview of the protocol conversions. Dashed lines denote a conversion via a third protocol.

376 the conversions. Some of them are implemented as a composition of two other protocols, e.g., for
 377 $A \rightarrow B$ we take the path $A \rightarrow Y \rightarrow B$.

378 A.3 Building Blocks for Neural Networks

379 Neural networks often have a very regular structure. The input is a tensor, a multi-dimensional
 380 array, which gets transformed by a sequence of different layers. Common are linear and pooling
 381 layers combined with (non-linear) activation functions. [19] We use the generic protocols (§ A.2) to
 382 construct building blocks to securely evaluate these layers on shared inputs. Since some are more
 383 efficiently computed with an arithmetic protocol and others with a Boolean protocol, we use the
 384 conversions (§ A.2) to change the representation of the data as necessary.

385 **Fixed-Point Arithmetic** Parameters, inputs, and outputs of neural networks are usually represented
 386 by floating-point numbers. While Boolean MPC can be used to evaluate Boolean circuits encoding
 387 the floating-point operations [12], the native operations of arithmetic MPC protocols are usually
 388 more efficient. Therefore, we encode decimal numbers into elements of \mathbb{Z}_{2^ℓ} using a fixed-point
 389 representation and the truncation protocol by [36] for A -shared values. It involves only local
 390 computation which makes it very efficient. The downside is that this protocol does not achieve perfect
 391 correctness and errors can be introduced.

392 **Linear Layers** For fully-connected, and convolution layers – both are essentially matrix multipli-
 393 cations, we use the arithmetic sharings A and α . While we can build matrix multiplication circuits
 394 out of additions and multiplications, it is more efficient to operate directly on (element-wise shared)
 395 matrices. The multiplication in both sharings generalize naturally to the multiplication of (element-
 396 wise shared) matrices: In GMW (A), we use multiplication triples $(\langle \mathbf{A} \rangle^A, \langle \mathbf{B} \rangle^A, \langle \mathbf{C} \rangle^A)$ where each
 397 element is a matrix. Then secure matrix multiplication works exactly as integer multiplication, but
 398 using element-wise addition and matrix multiplication instead of integer addition and multiplication,
 399 respectively [36]. This way, we mask each entry only once, instead of doing it separately for each
 400 integer multiplication, which saves communication. In the same fashion, the integer multiplication in
 401 the α sharing, can be generalized to matrix multiplication [38]. Then, $\Delta_{\mathbf{X}}, \delta_{\mathbf{X}}$ etc. are also matrices
 402 of appropriate size.

403 **ReLU** The ReLU operation ($\text{ReLU}(x) = \max(0, x)$) is a commonly used activation function.
 404 There are different approaches to compute a ReLU layer. First, note that, since the most significant bit
 405 $\text{msb}(x)$ encodes the sign of a number x in twos-complement, one can write $\text{ReLU}(x) = \neg \text{msb}(x) \cdot x$.
 406 Thus, the problem can be reduced to obtaining $\text{msb}(x)$ and computing the product. The MSB can be

407 computed either by converting the whole share into a Boolean sharing or via specialized bit extraction
408 protocols (e.g. [11]).

409 **MaxPool** To evaluate a MaxPool layer, we use a Boolean sharing (Y , B , or β) because of the
410 required comparisons. For each position of the window, the maximum is computed as follows: The
411 circuit consists of a balanced binary tree where each node is comprised of a $>$ comparison circuit
412 connected to a multiplexer such that the maximum of the input values is forwarded towards the root
413 of the tree. For the comparison, either size-optimized (for Y) or depth-optimized (for B, β) circuits
414 are used.

415 **AvgPool** AveragePool is a linear operation, consisting of a summation and an element-wise division
416 through the window (or kernel) size k . Since the window size k is part of the functionality, and, thus,
417 publicly known, we can write the latter as a multiplication with the constant $1/k$. However, since
418 $1/k < 1$ for all non-trivial cases, we need to take the fixed-point representation into account, and
419 execute a truncation protocol.

420 **Appendix B Extending the MOTION Framework**

421 This work builds upon and extends the MOTION framework for mixed-protocol multi-party com-
422 putation [7]. For more information about the original architecture of the framework, the reader is
423 referred to [7, § 4]. In the following we give a short overview of our extensions.

424 **Communication** The communication subsystem of the MOTION framework was *redesigned* to
425 allow easy-to-use and flexible asynchronous message-based communication. The challenge is that
426 multiple messages might concurrently be sent and expected on the sender and the receiver side,
427 respectively. So the system needs to make sure that every message ends up in the right place without
428 being able to rely on an order among the messages. Also – in the multi-party setting – messages
429 may be sent to and received from multiple parties. A low-level *transport* implements the actual
430 sending and receiving of messages between two parties, e.g., via TCP. The details of connection
431 setup and use as well as the used libraries are completely hidden from the user. Thus, the framework
432 can be easily adapted to use other transports instead, e.g., WebSockets or QUIC. The high-level
433 *CommunicationLayer* offers a simple but flexible API for sending messages to other parties. The
434 methods return immediately while the sending itself happens in the background. Received messages
435 are passed to message handlers which define how incoming messages of certain types are processed.
436 For synchronization between all parties, we provide a builtin *barrier* as synchronization mechanism.

437 **Protocol Implementation** Based on the terminology of circuits, two of the main concepts in the
438 MOTION framework [7] are *gates* and *wires*. Both are abstract interfaces which have been redesigned
439 to provide a cleaner API and reduced memory overhead. Wires are the passive components and hold
440 the local share of a secret-shared value. They can be seen as low-level variables and have a builtin
441 synchronization mechanism allowing a consumer to safely wait for it to obtain its value. A gate object
442 represents the active part of the computation encapsulating the protocol for a single operation, e.g.,
443 a primitive operation in a circuit or a tensor operation in a neural network. We introduce so-called
444 *protocol providers*, which bundle all the required functionality for a 2PC protocol with a common
445 interface. Hence, circuits can be constructed in a completely generic way: Given wires and an
446 identifier of an operation, a protocol provider will construct the corresponding gate and return the
447 output wire.

448 **Backend** MOTION [7] gathers all required components in so-called backend classes. The original
449 Backend in MOTION [7] contained a lot of unrelated functionality and was tightly coupled with the
450 rest of the framework. Now, the responsibility of the new backend classes has been reduced such that
451 they primarily construct and coordinate the required components (e.g., provider of OTs and protocols)
452 for their use-cases. Depending on the setting, different implementations for certain protocol can be

453 chosen. Currently, there are backends for the two-party and the multi-party settings, as well as a
454 specialized backend for evaluating tensor operations in neural networks.

455 **Execution** The encapsulation of the primitive operation in gate objects allows the decoupling of
456 the protocols description and its evaluation strategy. MOTION [7] uses *fibers* to evaluate gates, i.e.,
457 threads that run in user space opposed to standard threads managed by the operating system kernel.
458 They need less resources, and switching between fibers is done completely in user space. This offers
459 additional flexibility, e.g., allows using custom schedulers and allocators, and can yield improved
460 performance compared to threads [28]. We extracted the execution code into *executor* classes, that
461 takes a collection of gates and executes them according to some strategy, so that different strategies
462 can be implemented in different executors. As in [7], the executor for general-purpose MPC with
463 arbitrary circuits creates fibers for all gates which are then executed by a thread pool. For neural
464 networks, we use a different approach: Their computation graphs are often very simple or even
465 straight-line, and the tensor operations are more complex than the primitive operations in circuits.
466 Hence, we evaluate the gates sequentially while parallelizing the gate evaluations themselves.

467 **Statistics** MOTION [7] records detailed run-times during the execution, and computes the mean,
468 median, and standard deviation for repeated experiments. We added support to output this data
469 in the JSON format together with communication statistics and metadata of the experiments (e.g.,
470 experiment name, hostname, command line arguments, etc.). This makes it easy to import the data in
471 other software for further analysis.

472 **Implemented Protocols** MOTION [7] was released with three protocols for generic MPC with an
473 arbitrary number of parties. This work *additionally* implements support for the five 2PC protocols
474 discussed in § A.2. The original OT implementation by [7] was revised and extended, e.g., with
475 vectorization for C-OT (cf. [36]), and the OT-based multiplication protocols (§ A.1) have been
476 implemented on top. In addition to the general-purpose protocols, specialized building blocks (§ 2,
477 § A.3) have been implemented for the most common operations in neural networks. Depending on the
478 concrete operations, they are implemented using either the arithmetic or the Boolean 2PC protocols.

479 **File Formats** MOTION [7] allows building applications using its C++ interface and can also import
480 circuits from various simple file formats used in the MPC community [44, 1, 13]. So far, HyCC [9]
481 circuits were only rudimentary supported via a modified version of the HyCC adapter for ABY by [9].
482 A new `HyCCAdapter` has been developed that uses HyCC’s `libcircuit` to convert the HyCC circuits
483 into the respective MOTION data structures while hiding this from the MOTION user. Moreover,
484 a new `OnnxAdapter` enables support for neural network descriptions provided in the ONNX file
485 format [37], an open standard for deep learning models. With the ONNX library⁴, the adapter can
486 parse the description and automatically construct the respective tensor operations in MOTION. Both
487 new adapters can be optionally enabled at compile time of the MOTION framework if the respective
488 libraries are available.

489 Appendix C Performance Evaluation

490 C.1 Generic 2PC

491 **Garbling Engine** We measure the raw garbling and evaluation performance of our implementation
492 of the half-gate [48] garbling scheme with fixed-key AES [6, 20] and AES-NI, and benchmark it also
493 with various circuits of different complexity such as AES-128, SHA-256, ReLU, and comparisons.
494 It achieves a comparable performance to the EMP-Toolkit [45], and is $1.8\times$ faster when SIMD
495 operations are enabled. The garbling rate of a single thread is sufficient to saturate a 10 Gbit/s
496 network links.

⁴ONNX: <https://github.com/onnx/onnx>

497 **Boolean and Hybrid Circuits** We use AES-128 and SHA-256, to benchmark the performance of
498 the Boolean protocols B , β , and Y . Moreover, we use the hybrid circuits for biometric matching
499 generated by HyCC [9], which are evaluated using a combination of an arithmetic and a Boolean
500 protocol. In both cases, we compare with the ABY framework [13] and the generic N -party protocols
501 implemented in MOTION [7]. The overall performances are similar. Without SIMD, ABY [13]
502 always has the best run-times, which can be explained by the worse parallelization used in the
503 MOTION framework [7] (see below). With SIMD the performance of MOTION (including this work)
504 becomes competitive (up to $20\times$ better than ABY during the online phase in the LAN). SIMD is not
505 supported by the ABY backend for HyCC [9].

506 **Multi-Threading in MOTION** MOTION [7] uses fibers to evaluate the gates of a circuit. However,
507 it had not yet been examined how efficient this approach to parallelization is. We benchmark this
508 by evaluating the AES-128 circuit with the Boolean protocols. Regardless of the number of threads
509 the speedup stays below 2 for Yao’s protocol (Y). For the B and β protocols, we cannot see any
510 significant speedup. This is likely because Yao’s protocol uses cryptographic operations, whereas
511 the other protocols use only cheaper bit-wise operations in the gates. We conclude that fibers are a
512 good and easy-to-use method for implementing more complex protocols that need more computation
513 and (rounds of) interaction. They are also helpful for prototyping new protocols, but they are *not*
514 an efficient parallelization scheme for circuits built of cheap primitive operations. In such cases,
515 grouping the gates in layers and evaluating these multi-threaded is more advantageous (cf. [13]).

516 C.2 Neural Networks

517 Tables 1 and 2 contain our benchmark results for the CryptoNets [16] and MiniONN [33] neural
518 networks for the MNIST dataset [30]. We compare our results with ABY [13] and MOTION [7]. In
519 Table 3, we compare the online run-times for the MiniONN [33] neural network for the CIFAR-10
520 dataset [29] with GAZELLE [24], DELPHI [35], and CrypTFlow2 [41].

Table 1: Run-times in ms for evaluating the CryptoNets [16] neural network with ReLU using numbers of bitlength $\ell = 32$. This work with NN Ops uses specialized neural network operations (§ A.3), whereas all other categories use a hybrid circuit generated using HyCC [9]. With SIMD, 16 copies of the network are evaluated in parallel. All protocols were executed with $N = 2$ parties, and the best run-times are marked in bold.

Implementation	Protocol	LAN		WAN	
		Setup	Online	Setup	Online
ABY [13]	$A + B$	396.3	123.0	2 613.6	607.6
	$A + Y$	393.1	128.6	2 450.1	572.7
MOTION [7]	$A + B$	1 970.9	1 457.8	4 785.8	2 173.1
	$A + Y$	2 017.3	1 453.6	4 791.4	2 211.9
MOTION [7] w/ SIMD	$A + B$	916.9	101.4	2 889.7	279.8
	$A + Y$	897.2	96.8	2 935.9	248.7
this work	$A + B$	2 221.6	2 075.4	5 961.3	3 059.0
	$A + Y$	1 962.0	2 177.4	5 474.3	2 805.8
	$\alpha + \beta$	3 626.0	3 241.0	7 257.0	3 226.3
	$\alpha + Y$	3 525.6	3 261.0	7 400.3	3 371.6
this work w/ SIMD	$A + B$	650.9	145.3	2 810.5	263.3
	$A + Y$	645.5	142.2	2 777.3	223.7
	$\alpha + \beta$	572.8	202.3	2 807.6	251.9
	$\alpha + Y$	590.1	210.6	2 667.3	254.4
this work w/ NN Ops.	$A + B$	111.2	11.1	1 005.7	1 047.4
	$A + Y$	103.4	11.1	1 062.4	858.6
	$\alpha + \beta$	129.1	28.1	1 766.7	982.2
	$\alpha + Y$	136.6	8.7	1 539.4	858.6

Table 2: Run-times in ms for evaluating the MiniONN [33] MNIST neural network using numbers of bitlength $\ell = 32$. This work with NN Ops uses specialized neural network operations (§ A.3), whereas all other categories use a hybrid circuit generated using HyCC [9]. With SIMD, 16 copies of the network are evaluated in parallel. All protocols were executed with $N = 2$ parties, and the best run-times are marked in bold.

Implementation	Protocol	LAN		WAN	
		Setup	Online	Setup	Online
ABY [13]	$A + B$	1 516.3	2 997.2	4 903.5	8 164.1
	$A + Y$	1 559.0	1 453.5	4 960.0	4 194.1
this work	$A + B$	6 423.3	49 712.7	10 391.8	46 298.3
	$A + Y$	16 138.6	15 073.4	19 085.8	16 387.4
	$\alpha + \beta$	109 778.9	48 342.2	110 201.4	48 562.7
	$\alpha + Y$	15 939.9	14 726.1	19 235.7	15 867.0
this work w/ SIMD	$A + B$	1 239.6	3 037.3	5 163.2	2 903.6
	$A + Y$	1 514.1	936.8	6 129.0	1 467.3
	$\alpha + \beta$	7 592.7	3 324.8	9 330.1	2 954.8
	$\alpha + Y$	1 315.6	928.6	6 664.4	973.1
this work w/ NN Ops.	$A + B$	649.7	171.5	4 826.6	5 138.1
	$A + Y$	523.8	126.7	4 476.5	2 574.9
	$\alpha + \beta$	737.0	151.6	6 424.5	4 654.0
	$\alpha + Y$	676.1	81.8	5 908.9	1 235.7

Table 3: Online Run-times in ms for evaluating MiniONN [33] CIFAR-10 neural network with ℓ -bit numbers. The run-time of GAZELLE is taken from the corresponding publication [24], and the numbers for DELPHI [35] and CrypTFlow2 are taken from the CrypTFlow2 paper [41]. Note that they were obtained in a different experimental environment. All protocols were executed with $N = 2$ parties. The best run-times are marked in bold, although different bit sizes ℓ are used by the different implementations.

Implementation	Protocol	ℓ	LAN	WAN
GAZELLE [24]	$A + Y$	60	3 560.0	—
DELPHI [35]	$A + Y$	41	$\approx 4\,070.0$	—
CrypTFlow2 [41]	$A + OT$	41	\approx 420.0	—
this work w/ NN Ops	$A + B$	32	1 167.7	20 773.1
	$A + Y$	32	1 246.8	34 783.5
	$\alpha + \beta$	32	624.2	7 882.5
	$\alpha + Y$	32	695.6	6 756.2
	$A + B$	64	2 024.6	37 584.9
	$A + Y$	64	2 231.1	66 187.8
	$\alpha + \beta$	64	989.5	11 768.9
	$\alpha + Y$	64	1 226.6	10 591.8