

---

# CRYPTEN: Secure Multi-Party Computation Meets Machine Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Secure multi-party computation (MPC) allows parties to perform computations  
2 on data while keeping that data private. This capability has great potential for  
3 machine-learning applications: it facilitates training of machine-learning models  
4 on private data sets owned by different parties, evaluation of one party’s private  
5 model using another party’s private data, *etc.* Although a range of studies imple-  
6 ment machine-learning models via secure MPC, such implementations are not  
7 yet mainstream. Adoption of secure MPC is hampered by the absence of flexible  
8 software frameworks that “speak the language” of machine-learning researchers  
9 and engineers. To foster adoption of secure MPC in machine learning, we present  
10 CRYPTEN: a software framework that exposes popular secure MPC primitives via  
11 abstractions that are common in modern machine-learning frameworks, such as  
12 tensor computations, automatic differentiation, and modular neural networks. This  
13 paper describes the design of CRYPTEN and measure its performance on state-of-  
14 the-art models for text classification, speech recognition, and image classification.  
15 Our benchmarks show that CRYPTEN’s GPU support and high-performance com-  
16 munication between (an arbitrary number of) parties allows it to perform efficient  
17 private evaluation of modern machine-learning models under a *semi-honest* threat  
18 model. For example, two parties using CRYPTEN can securely predict phonemes  
19 in speech recordings using Wav2Letter [17] faster than real-time. We hope that  
20 CRYPTEN will spur adoption of secure MPC in the machine-learning community.

## 21 1 Introduction

22 Secure multi-party computation (MPC; [29, 68]) allows parties to collaboratively perform computa-  
23 tions on their combined data sets without revealing the data they possess to each other. This capability  
24 of secure MPC has the potential to unlock a variety of machine-learning applications that are currently  
25 infeasible because of data privacy concerns. For example, secure MPC can allow medical research  
26 institutions to jointly train better diagnostic models without having to share their sensitive patient  
27 data [26] or allow social scientists to analyze gender wage gap statistics without companies having  
28 to share sensitive salary data [41]. The prospect of such applications of machine learning with  
29 rigorous privacy and security guarantees has spurred a number of studies on machine learning via  
30 secure MPC [37, 40, 47, 57, 62, 65, 66]. However, at present, adoption of secure MPC in machine  
31 learning is still relatively limited considering its wide-ranging potential. One of the main obstacles to  
32 widespread adoption is that the complexity of secure MPC techniques puts them out of reach for most  
33 machine-learning researchers, who frequently lack in-depth knowledge of cryptographic techniques.

34 To foster the adoption of secure MPC techniques in machine learning, we present CRYPTEN: a  
35 flexible software framework that aims to make modern secure MPC techniques accessible to machine-  
36 learning researchers and developers without a background in cryptography. Specifically, CRYPTEN

37 provides a comprehensive tensor-computation library in which all computations are performed via  
38 secure MPC. CRYPTEN’s API closely follows the API of the popular PyTorch framework for machine  
39 learning [53, 54], which makes it easy to use for machine-learning practitioners. For example,  
40 it provides automatic differentiation and a modular neural-network package. CRYPTEN assumes  
41 an *semi-honest* threat model [29, §2.3.2] and works for an arbitrary number of parties. To make  
42 private training and inference efficient, CRYPTEN off-loads computations to the GPU and uses  
43 high-performance communication libraries to implement interactions between parties.

44 The paper presents: (1) an overview of CRYPTEN’s design principles; (2) a description of the  
45 design of CRYPTEN and of the secure MPC protocols implemented; (3) a collection of benchmark  
46 experiments using CRYPTEN to run private versions of state-of-the-art models for text classification,  
47 speech recognition, and image classification; and (4) a discussion of open problems and a roadmap  
48 for the further development of CRYPTEN. Altogether, the paper demonstrates that CRYPTEN’s  
49 flexible, PyTorch-like API makes private inference and training of modern machine-learning models  
50 easy to implement and efficient. For example, CRYPTEN allows two parties to privately classify  
51 an image [25, 34] in 2-3 seconds, or to securely make phoneme predictions for 16kHz speech  
52 recordings [17] faster than real-time. We hope that CRYPTEN’s promising performance and ease-of-  
53 use will foster the adoption of secure MPC by the machine-learning community, and pave the way  
54 for a new generation of secure and private machine-learning systems.

## 55 2 Related Work

56 This work is part of a larger body of work that develops secure MPC systems for machine learning.  
57 Most closely related to our work is CryptGPU [62], which implements an 2-out-of-3 replicated secret  
58 sharing protocol [4, 36] *on top of* CRYPTEN. This protocol provides security against semi-honest  
59 corruption, but also limits it to the three-party setting. CryptGPU is part of a larger family of three-  
60 party protocols that provides malicious security. For example, Falcon [66] implements a maliciously  
61 secure MPC protocol for the three-party setting, combining techniques from SecureNN [65] and  
62 ABY3 [47]. Falcon allows evaluation and training of convolutional networks such as AlexNet [39] and  
63 VGG [61]. Other systems that work in this setting include Astra [16], Blaze [55], and CryptFlow [40].

64 There also exists a family of two-party systems that, like CRYPTEN, assume a *semi-honest* threat  
65 model. These systems include Gazelle [37], Chameleon [57], EzPC [15], MiniONN [44], Se-  
66 cureML [48], PySyft [59], and Delphi [46]. XONN [58] also works in the two-party setting but  
67 provides malicious security. Compared to these systems, CRYPTEN provides a more flexible machine-  
68 learning focused API<sup>1</sup> that supports reverse-mode automatic differentiation, implements a rich set of  
69 functions, and natively runs on GPUs. Moreover, CRYPTEN supports a wider range of use cases by  
70 working with an arbitrary number of parties, and make communication between parties efficient via  
71 communication primitives that were optimized for high-performance distributed computing.

## 72 3 Design Principles

73 In the development of CRYPTEN, we adopted the following two main design principles:

74 **Machine-learning first API.** CRYPTEN has a general purpose, machine-learning first API design.  
75 Most other secure MPC frameworks [33] adopt an API that stays close to the underlying MPC  
76 protocols. This hampers adoption of these frameworks in machine learning, for example, because  
77 they do not natively support tensor operations (but only scalar operations) and because they lack  
78 features that machine-learning researchers have come to expect, such as automatic differentiation.  
79 Instead, CRYPTEN implements the tensor-computation API of the popular PyTorch machine-learning  
80 framework [53], implements reverse-mode automatic differentiation, provides a modular neural-  
81 network package with corresponding learning routines, and supports GPU computations. We aim to  
82 allow developers to transition code from PyTorch to CRYPTEN by changing a single Python `import`.

83 **Eager execution.** CRYPTEN adopts an imperative programming model. This is different from  
84 existing MPC frameworks, which generally implement compilers for their own domain-specific  
85 languages [33]. While compiler approaches have potential performance benefits, they slow down the

---

<sup>1</sup>CryptFlow [40] also provides such an API by integrating deeply with TensorFlow [1], but unlike CRYPTEN,  
it does not support PyTorch’s eager execution model [54] or GPU support.

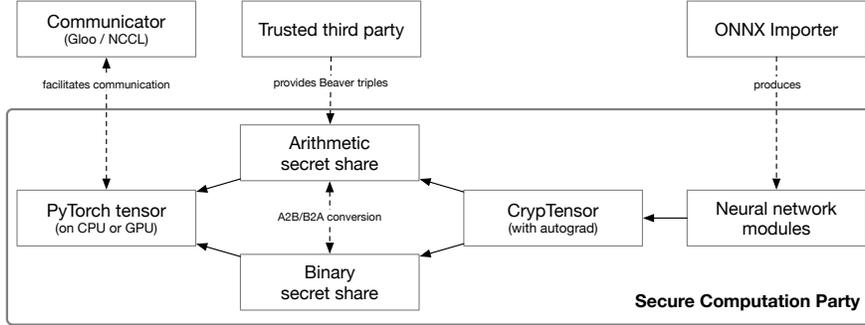


Figure 1: High-level overview of the design of CRYPTEN. See text in Section 4 for details.

86 development cycle, make debugging harder, and prevent users from using arbitrary host-language  
 87 constructs [3]. Instead, CRYPTEN follows the recent trend in machine learning away from graph  
 88 compilers [1] to frameworks that eagerly execute computations [3, 54], providing a better developer  
 89 experience. Yet, CRYPTEN is performant because it implements state-of-the-art secure MPC protocols  
 90 (for settings with arbitrary number of parties), because it uses PyTorch’s highly optimized tensor  
 91 library for most computations, because computations can be off-loaded to the GPU, and because it  
 92 uses communication libraries that were optimized for high-performance distributed computing.

## 93 4 Design Overview

94 Figure 1 gives an overview of CRYPTEN’s design. Parties perform computations using efficient  
 95 PyTorch tensor operations. Because secure MPC computations are integer computations that are not  
 96 natively supported on GPUs, CRYPTEN maps between integer and floating-point computations on  
 97 GPUs; see Section 5.3. The multi-party computations are implemented on arithmetic and binary  
 98 secret shares [22, 31]; see Section 5.1. Whereas many computations can be performed directly  
 99 on arithmetic secret shares, others require conversion between arithmetic and binary secret shares  
 100 (A2B) and back (B2A); see Section 5.2. Some multi-party computations require interaction  
 101 between parties via a *communicator* that employs the high-performance communication primitives  
 102 in Gloo [30] and NCCL [50]. Some multi-party computations require Beaver triples [7], which  
 103 are supplied by a *trusted third party* (TTP).<sup>2</sup>

114 All secure computations are wrapped in a `CrypTensor` object that implements the PyTorch tensor  
 115 API and that provides reverse-mode automatic differentiation (autograd) to enable gradient-based  
 116 training of arbitrary (deep) learning models. Figure 2 illustrates `CrypTensor` creation, *i.e.*, how  
 117 tensors are secret-shared and revealed, as well as a simple computation (addition). Note that each  
 118 party involved in the multi-party computation executes the same code. Whenever communication  
 119 between the parties is required (*e.g.*, as part of private multiplications), the communication acts  
 120 as a synchronization point between the parties. The `crypten.init()` call is required once to  
 121 establish the communication channel. In the example, the input tensor for the creation of the  
 122 arithmetic secret share is provided party `src=0`, which indicates the rank<sup>3</sup> of the party that  
 123 supplies the data to be secret-shared (the other parties executing this code may provide `None` as  
 input).

---

```

import crypten, torch

# set up communication and sync random seeds:
crypten.init()

# secret share tensor:
x = torch.tensor([1.0, 2.0, 3.0])
x_enc = crypten.cryptensor(x, src=0)

# reveal secret shared tensor:
x_dec = x_enc.get_plain_text()
assert torch.all_close(x_dec, x)

# add secret shared tensors:
y = torch.tensor([2.0, 3.0, 4.0])
y_enc = crypten.cryptensor(y, src=0)
xy_enc = x_enc + y_enc
xy_dec = xy_enc.get_plain_text()
assert torch.all_close(xy_dec, x + y)

```

---

Figure 2: Example of secret-sharing tensors, revealing tensors, and private addition in CRYPTEN.

<sup>2</sup>CRYPTEN adopts a trusted third party for generating Beaver triples for efficiency reasons, but we are planning to add TTP-free solutions based on additive homomorphic encryption [51] or oblivious transfer [38].

<sup>3</sup>CRYPTEN relies on MPI primitives for communication: each party knows their rank and the world size.

124 To enable deep-learning use cases, CRYPTEN  
 125 allows implementing neural networks following  
 126 PyTorch’s API. Figure 3 shows how to create  
 127 and encrypt neural networks and how to use  
 128 automatic differentiation in CRYPTEN. The ex-  
 129 ample assumes that some training sample and  
 130 the associated target label are provided by the  
 131 party with rank 0 (note the value of src). As  
 132 illustrated by the example, CRYPTEN’s API  
 133 closely follows that of PyTorch. Indeed, it is  
 134 possible to write a single training loop that can  
 135 be used to train models using CRYPTEN or Py-  
 136 Torch without code changes. This makes it easy  
 137 to adapt PyTorch code to use secure MPC for  
 138 its computations, and it also makes debugging  
 139 easier. The appendix presents a table listing all  
 140 tensor functions that CRYPTEN implements.

141 To enable interoperability with existing machine-  
 142 learning platforms, neural networks can be im-  
 143 ported into CRYPTEN via ONNX. Figure 4  
 144 shows how a PyTorch model is imported  
 145 into CRYPTEN. The example illustrates how  
 146 CRYPTEN makes private inference with a  
 147 ResNet-18 easy. The example in the figure also  
 148 demonstrates CRYPTEN’s GPU support. One  
 149 caveat is that all parties must use the same type  
 150 of device (*i.e.*, CPU or GPU) for computations.

## 151 5 Secure Computations

152 To facilitate secure computations, CRYPTEN  
 153 implements arithmetic secret sharing [22] and  
 154 binary secret sharing [31], as well as conver-  
 155 sions between these two types of sharing [23].  
 156 Arithmetic secret sharing is particularly well-  
 157 suited for operations that are common in mod-  
 158 ern machine-learning models, such as matrix  
 159 multiplications and convolutions. Binary se-  
 160 cret sharing is required for evaluating certain  
 161 other common functions, such as rectified lin-  
 162 ear units. We provide a high-level overview of  
 163 CRYPTEN’s secure computation protocol here; a  
 164 detailed description is presented in the appendix.

### 165 5.1 Secret Sharing

166 **Arithmetic secret sharing** shares a scalar value  $x \in \mathbb{Z}/Q\mathbb{Z}$ , where  $\mathbb{Z}/Q\mathbb{Z}$  denotes a ring with  $Q$   
 167 elements, across parties  $p \in \mathcal{P}$ . We denote the sharing of  $x$  by  $\langle x \rangle = \{[x]_p\}_{p \in \mathcal{P}}$ , where  $[x]_p \in \mathbb{Z}/Q\mathbb{Z}$   
 168 indicates party  $p$ ’s share of  $x$ . The shares are constructed such that their sum reconstructs the original  
 169 value  $x$ , that is,  $x = \sum_{p \in \mathcal{P}} [x]_p \pmod{Q}$ . To share a value  $x$ , the parties generate a pseudorandom  
 170 zero-share [18] with  $|\mathcal{P}|$  random numbers that sum to 0. The party that possesses the value  $x$  adds  $x$  to  
 171 their share and discards  $x$ . We use a fixed-point encoding to obtain  $x$  from a floating-point value,  $x_R$ .  
 172 To do so, we multiply  $x_R$  with a large scaling factor  $B$  and round to the nearest integer:  $x = \lfloor Bx_R \rfloor$ ,  
 173 where  $B = 2^L$  for some precision of  $L$  bits. To decode a value,  $x$ , we compute  $x_R \approx x/B$ .

174 **Binary secret sharing** is a special case of arithmetic secret sharing that operates within the binary  
 175 field  $\mathbb{Z}/2\mathbb{Z}$ . A binary secret share,  $\langle x \rangle$ , of a value  $x$  is formed by arithmetic secret shares of the bits  
 176 of  $x$ , setting  $Q=2$ . Each party  $p \in \mathcal{P}$  holds a share,  $\langle x \rangle_p$ , such that  $x = \bigoplus_{p \in \mathcal{P}} \langle x \rangle_p$  is satisfied.

---

```

import crypten.optimizer as optimizer
import crypten.nn as nn

# create model, criterion, and optimizer:
model_enc = nn.Sequential(
    nn.Linear(sample_dim, hidden_dim),
    nn.ReLU(),
    nn.Linear(hidden_dim, num_classes),
).encrypt()
criterion = nn.CrossEntropyLoss()
optimizer = optimizer.SGD(
    model_enc.parameters(), lr=0.1, momentum=0.9,
)

# perform prediction on sample:
target_enc = crypten.cryptensor(target, src=0)
sample_enc = crypten.cryptensor(sample, src=0)
output_enc = model_enc(sample_enc)

# perform backward pass and update parameters:
model_enc.zero_grad()
loss_enc = criterion(output_enc, target_enc)
loss_enc.backward()
optimizer.step()

```

---

Figure 3: Example using neural networks and automatic differentiation in CRYPTEN.

---

```

import torchvision.datasets as datasets
import torchvision.models as models
import torchvision.transforms as transforms

# download and set up ImageNet dataset:
transform = transforms.ToTensor()
dataset = datasets.ImageNet(
    imagenet_folder, transform=transform,
)

# secret share pre-trained ResNet-18 on GPU:
model = models.resnet18(pretrained=True)
model_enc = crypten.nn.from_pytorch(
    model, dataset[0],
).encrypt().cuda()

# perform inference on secret-shared images:
for image in dataset:
    image_enc = crypten.cryptensor(image).cuda()
    output_enc = model_enc(image_enc)
    output = output_enc.get_plain_text()

```

---

Figure 4: Private inference on secret-shared images using a secret-shared ResNet-18 model on GPU.

177 **Conversion from  $[x]$  to  $\langle x \rangle$**  is implemented by having the parties create a binary secret share of  
 178 their  $[x]_p$  shares, and summing the resulting binary shares. Specifically, the parties create a binary  
 179 secret share,  $\langle [x]_p \rangle$ , of all the bits in  $[x]_p$ . Subsequently, the parties compute  $\langle x \rangle = \sum_{p \in \mathcal{P}} \langle [x]_p \rangle$   
 180 using a Ripple-carry adder in  $\log_2(|\mathcal{P}|) \log_2(L)$  communication rounds [14, 21].

181 **Conversion from  $\langle x \rangle$  to  $[x]$**  is achieved by computing  $[x] = \sum_{b=1}^B 2^b [\langle x \rangle^{(b)}]$ , where  $\langle x \rangle^{(b)}$  denotes  
 182 the  $b$ -th bit of the binary share  $\langle x \rangle$  and  $B$  is the total number of bits in the shared secret,  $\langle x \rangle$ . To  
 183 create an arithmetic share of a bit, the parties use secret shares,  $([r^{(b)}], \langle r^{(b)} \rangle)$ , of random bits  $r^{(b)}$ .  
 184 The random bits are provided by the TTP, but we plan to add an implementation that generates them  
 185 off-line via oblivious transfer [38]. The parties use  $\langle r^{(b)} \rangle$  to mask  $\langle x \rangle^{(b)}$  and reveal the resulting  
 186 masked bit  $z^{(b)}$ . Subsequently, they compute  $[\langle x \rangle^{(b)}] = [r^{(b)}] + z^{(b)} - 2[r^{(b)}]z^{(b)}$ .

## 187 5.2 Secure Computation

188 Arithmetic and binary secret shares have homomorphic properties that can be used to implement  
 189 secure computations. All computations in CRYPTEN are based on private addition and multiplication.

190 **Private addition** of two arithmetically secret shared values,  $[z] = [x] + [y]$ , is implemented by  
 191 having each party  $p$  sum their shares of  $[x]$  and  $[y]$ : each party  $p \in \mathcal{P}$  computes  $[z]_p = [x]_p + [y]_p$ .

192 **Private multiplication** is implemented using random Beaver triples [7],  $([a], [b], [c])$  with  $c = ab$ , that  
 193 are provided by the TTP. The parties compute  $[\epsilon] = [x] - [a]$  and  $[\delta] = [y] - [b]$ , and decrypt  $\epsilon$  and  $\delta$   
 194 without information leakage due to the masking. They compute the result  $[x][y] = [c] + \epsilon[b] + [a]\delta + \epsilon\delta$ ,  
 195 using trivial implementations of addition and multiplication of secret shares with public values.

196 **Linear functions** are trivially implemented as combinations of private addition and multiplication.  
 197 This allows CRYPTEN to compute dot products, outer products, matrix products, and convolutions.

198 **Non-linear functions** are implemented using standard approximations that only require private addi-  
 199 tion and multiplication. Specifically, CRYPTEN evaluates exponentials using a limit approximation,  
 200 logarithms using Householder iterations [35], and reciprocals using Newton-Rhapson iterations.  
 201 This allows CRYPTEN to implement functions that are commonly used in machine-learning models,  
 202 including the sigmoid, softmax, and logistic-loss functions, as well as their gradients.

203 **Comparators** are implemented using a function that evaluates  $[z < 0]$  by: (1) converting  $[z]$  to a  
 204 binary secret-share  $\langle z \rangle$ ; (2) computing its sign bit,  $\langle b \rangle = \langle z \rangle \gg (L - 1)$ ; and (3) converting the  
 205 resulting bit to an arithmetic sharing  $[b]$ . This function allows CRYPTEN to implement arbitrary  
 206 comparators. For example, it evaluates  $[x < y]$  by computing  $[z] = [x] - [y]$  and evaluating  $[z < 0]$ .  
 207 Similarly, CRYPTEN can evaluate: (1) the sign function via  $\text{sign}([x]) = 2[x > 0] - 1$ ; (2) the absolute  
 208 value function via  $|[x]| = [x] \text{sign}([x])$ ; and (3) rectified linear units via  $\text{ReLU}([x]) = [x][x > 0]$ .  
 209 CRYPTEN also supports multiplexing; to do so, it evaluates  $[c ? x : y] = [c][x] + (1 - [c])[y]$ .

210 **Lemma 1.** *The CRYPTEN secure-computation protocol is secure against information leakage against*  
 211 *any static passive adversary corrupting up to  $|\mathcal{P}| - 1$  of the  $|\mathcal{P}|$  parties involved in the computation.*

212 The proof of this lemma follows trivially from [9, 11, 21, 23], and is given in the appendix. We adopt  
 213 a protocol that provides security under a *semi-honest* threat model because it enables a wide range of  
 214 use cases of secure machine learning, whilst being more efficient than maliciously secure protocols.

## 215 5.3 Off-loading Computations to the GPU

216 Hardware acceleration via GPUs is a critical component for training and inference in modern machine-  
 217 learning models. Akin to frameworks such as PyTorch [54] and TensorFlow [1], CRYPTEN can  
 218 off-load computations to the GPU. On the GPU, it uses highly-optimized implementations for a range  
 219 of functions that are provided by CUDA libraries such as cuBLAS [19] and cuDNN [20].

220 Unfortunately, these libraries are designed for computations on floating-point numbers and do not  
 221 support the integer types required to perform computations on  $L$ -bit fixed-point numbers. Akin  
 222 to [62], we circumvent this problem by observing that for all integers  $a, b \in \mathbb{Z} \cap [-2^{26}, 2^{26}]$ , we  
 223 can compute the product  $ab$  using 64-bit floating-point representations and still recover the correct  
 224 value over the integers. Specifically, CRYPTEN splits each 64-bit variable into four components,  
 225  $a = a_0 + 2^{16}a_1 + 2^{32}a_2 + 2^{48}a_3$ , where each  $a_i$  represents a 16-bit integer component. We  
 226 compute a product  $ab$  of 64-bit integers by summing 10 pairwise products of their 16-bit components.

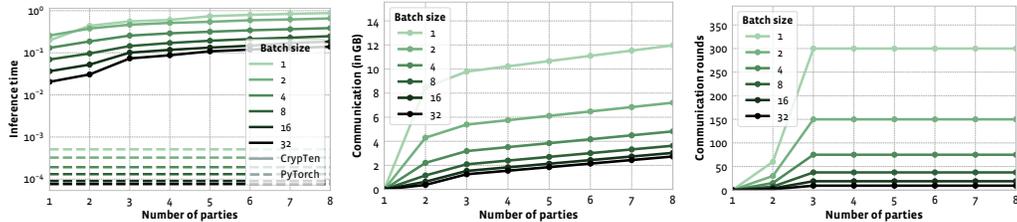


Figure 5: Benchmarks for inference with text-sentiment classification model on GPUs in CRYPTEN and PyTorch. **Left:** Average wall-clock time per sample (in seconds). **Middle:** Number of bytes communicated per sample, per party (in GB). **Right:** Number of communication rounds per sample.

227 The pairwise products of the 16-bit components are computed in parallel using highly optimized  
 228 floating-point CUDA kernels. The same approach is used for matrix multiplications and convolutions.  
 229 CRYPTEN further optimizes this approach by splitting into only 3 components of 22-bits each when  
 230 possible, which reduces the number of pairwise products required to 6 (see [62, Remark II.1]).

## 231 6 Benchmarks

232 To measure the performance of CRYPTEN, we performed experiments on three tasks: (1) text  
 233 classification using a linear model that learns word embeddings; (2) speech recognition using  
 234 the Wav2Letter model [17]; and (3) image classification using residual networks [34] and vision  
 235 transformers [25]. Because of space constraints, we focus on private inference using a secret-shared  
 236 model on secret-shared data here, but our benchmark results with private training are very similar.

237 We performed benchmark experiments on a proprietary cluster, testing inference on both CPUs (Intel  
 238 Skylake 18-core 1.6GHz) and GPUs (nVidia P100). We set the number of OpenMP threads to 1 in all  
 239 benchmarks. All experiments were performed with the parties running in separate processes on a  
 240 single machine. For GPU experiments, each party was assigned its own GPU. Although this setup is  
 241 faster than a scenario in which each party operates its own machine,<sup>4</sup> we believe our benchmark results  
 242 provide a good sense of CRYPTEN’s performance. We average computation times over 30 batches,  
 243 excluding the computation on the first batch as that computation may include CuDNN benchmarking.  
 244 Code reproducing the results of our experiments is available on <https://anonymized.ai>.

245 In our benchmarks, we focus on comparing (ciphertext) CRYPTEN computation with (plaintext)  
 246 PyTorch computation. We refer the reader to [32, 62] for benchmarks that compare CRYPTEN to  
 247 other secure MPC frameworks. Specifically, [32] finds CRYPTEN is 11-18× faster than PySyft [59]  
 248 and approximately 3× faster than TF-Trusted [13] in MNIST classification [42] on CPU.

### 249 6.1 Text Classification

250 We performed text-sentiment classification experiments on the Yelp review dataset [69] using a model  
 251 that consists of a linear layer operating on word embeddings. The embedding layer contains 32-  
 252 dimensional embeddings of 519, 820 words, and the linear layer produces a binary output indicating  
 253 the sentiment of the review. We evaluated the model on GPUs, varying the batch size and the number  
 254 of parties participating. The normalized mean squared error ( $\frac{\|x-y\|^2}{\|x\|^2}$ ) between the output of the  
 255 CRYPTEN model and that of its PyTorch counterpart was smaller than  $4 \cdot 10^{-4}$  in all experiments.

256 Figure 5 presents the results of our experiments. The figure shows inference time *per sample* (in  
 257 seconds) as a function of the number of parties involved in the computation for varying batch sizes  
 258 (left); the amount of communication required per sample, *per party* (in GB); and the number of  
 259 communication rounds required per sample. We include results in which the number of parties is 1:  
 260 herein, we run the CRYPTEN protocol but involve no other parties, which implies that the single party  
 261 is running the protocol on unencrypted data. One-party results allow us to bisect different sources of  
 262 computational overhead: specifically, they separate overhead due to communication from overhead  
 263 due to fixed-point encoding, function approximations, and (lack of) sparse-matrix operations.

<sup>4</sup>Communication between GPUs in two machines connected via InfiniBand has approximately 20× lower throughput than communication between two GPUs in the same machine via NVLink (25GB/s versus 600GB/s).

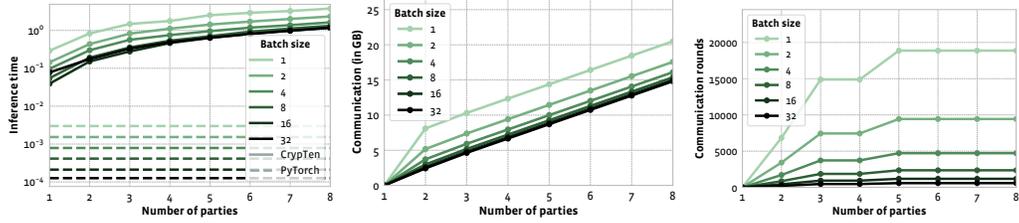


Figure 6: Benchmarks for inference with Wav2Letter model on GPUs in CRYPTEN and PyTorch. **Left:** Average wall-clock time per sample (in seconds). **Middle:** Number of bytes communicated per sample, per party (in GB). **Right:** Number of communication rounds per sample.

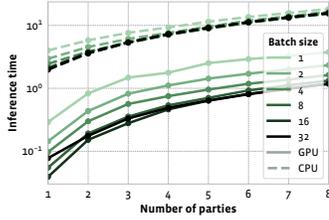


Figure 7: Wall-clock time per sample (in sec.) for Wav2Letter inference on CPUs and GPUs.

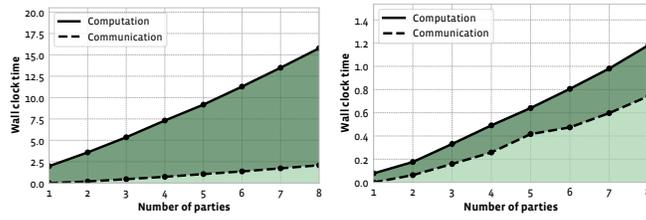


Figure 8: Average wall-clock time per sample (in seconds) for communication and computation during inference with Wav2Letter model on CPU (**left**) and GPU (**right**).

264 The results in Figure 5 show that CRYPTEN is about 2.5–3 orders of magnitude slower than PyTorch  
 265 in text-sentiment classification, depending on the number of parties involved. Most computational  
 266 overhead is the word embedding layer: whereas PyTorch can evaluate this layer efficiently via a  
 267 sparse matrix multiplication, CRYPTEN cannot do sparse lookups as they would reveal information  
 268 on the encrypted input. Instead, CRYPTEN performs a full matrix multiplication between the word-  
 269 count vector and the embedding matrix. Yet, text sentiment predictions are quite fast in CRYPTEN:  
 270 inference takes only 0.03 seconds per sample in the two-party setting with a batch size of 32.

271 The results also show that increasing the batch size is an effective way to reduce inference time  
 272 and communication per sample. The number of communication rounds is independent of the batch  
 273 size, which means communication rounds can be amortized by using larger batch sizes. The number  
 274 of bytes communicated is partly amortized as well because the size of weight tensors (*e.g.*, in  
 275 linear layers) does not depend on batch size. The results also show that whereas the number of  
 276 communication rounds increases when moving from two-party to three-party computation, it remains  
 277 constant afterwards. The larger number of communication rounds for three-party computation stems  
 278 from the public division protocol, which requires additional communication rounds when more than  
 279 two parties are involved to prevent wrap-around errors (see the appendix for details).

## 280 6.2 Speech Recognition

281 We performed speech-recognition experiments using Wav2Letter [17] on the LibriSpeech dataset [52].  
 282 The LibriSpeech dataset contains 16 kHz audio clips represented as a waveform (16,000 samples per  
 283 second). Because the audio clips vary in length, we clip all of them to 1 second for the benchmark.  
 284 Wav2Letter is a network with 13 convolutional layers using rectified linear unit (ReLU; [49]) activa-  
 285 tions.<sup>5</sup> The network operates directly on the waveform input, predicting one of 29 labels (26 letters  
 286 plus 3 special characters). The first two layers use a filter size of 250 (with stride 160) and 48 (stride  
 287 2). The next seven layers use filter size 7, followed by two layers with filter size 32 and 1 (all with  
 288 stride 1). All layers except the last two have 250 channels. The last two layers have 2,000 channels.

289 The results in Figure 6 show that CRYPTEN is about 2.5–3 orders of magnitude slower than PyTorch  
 290 depending on the number of parties involved. For Wav2Letter, the overhead is largely due to the  
 291 ReLU layers in the network: evaluating a ReLU function requires a comparison, which involves a

<sup>5</sup>We used the reference implementation of Wav2Letter in `torchaudio`.

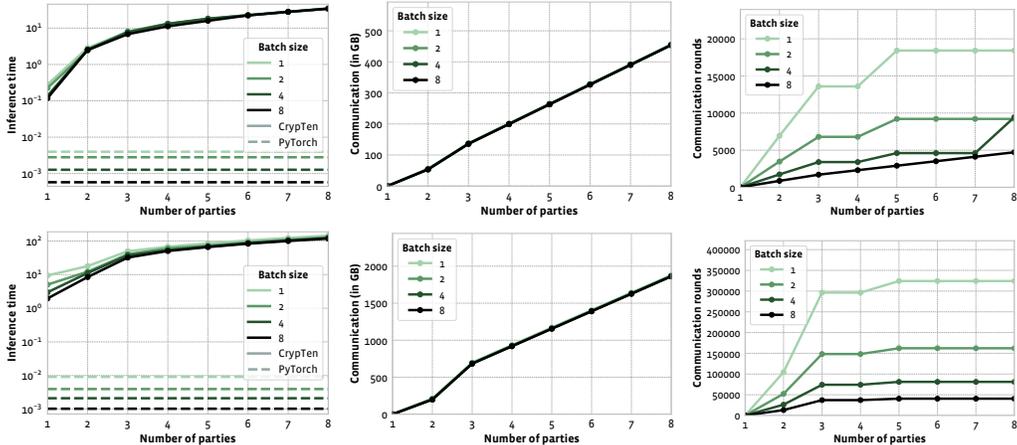


Figure 9: Benchmarks for inference with image-classification models on GPUs in CRYPTEN and PyTorch. **Top:** Results for ResNet-18 model. **Bottom:** Results for ViT-B/16 vision transformer. **Left:** Average wall-clock time per sample (in seconds). **Middle:** Number of bytes communicated per sample, per party (in GB). **Right:** Number of communication rounds per sample.

292 conversion between arithmetic and binary secret sharing and back (see the appendix). The number  
 293 of communication rounds increases when the number of parties grows beyond 4: CRYPTEN uses  
 294 a tree reduction for the summation in the comparator protocol, which implies that the number of  
 295 communication rounds grows whenever the number of parties increases from  $2^k$  to  $2^{k+1}$ .

296 Figure 7 also presents results comparing Wav2Letter inference time between CPUs and GPUs.  
 297 The results in the figure show that CRYPTEN is 1-2 orders of magnitude faster on GPUs than on  
 298 CPUs. In real-world settings, this speedup can make the difference between a secure MPC use case  
 299 being practical or not. Figure 8 shows how much wall-clock time is spent on communication and  
 300 computation, respectively, when performing inference with Wav2Letter (using batch size 32). The  
 301 results suggest that, whereas multi-party evaluation is compute-bound on CPU, it is communication-  
 302 bound on GPU. On GPUs, 63% of the time is spent on communication in eight-party computation.

### 303 6.3 Image Classification

304 We performed image-classification experiments on the ImageNet dataset using residual networks  
 305 (ResNets; [34]) and vision transformers (ViT; [25]).<sup>6</sup> We experimented with a ResNet-18 with 18  
 306 convolutional layers and with a ViT-B/16 model that has 12 multi-head self-attention layers with  
 307 12 heads each, operating on image patches of  $16 \times 16$  pixels. Following common practice [34], we  
 308 preprocess images by rescaling them to size  $256 \times 256$  and taking a center crop of size  $224 \times 224$ .

309 Figure 9 presents the results of our image-classification benchmarks, which show that two parties  
 310 can securely evaluate a ResNet-18 model in 2.49 seconds and a ViT-B/16 model in 8.47 seconds.  
 311 A notable difference compared to the prior results is that the number of bytes communicated per  
 312 sample is no longer reduced by increasing the batch size. The reason for this is that the vast majority  
 313 of communication involves tensors that have the same size as intermediate activation functions:  
 314 activation tensors are much larger than weight tensors in image-classification models. The amount  
 315 of communication required to evaluate the ViT-B/16 model is particularly high due to the repeated  
 316 evaluation of the softmax function in the attention layer of Transformers [64], which involves a limit  
 317 approximation for the exponential and a Newton-Rhapson approximation for the normalization. We  
 318 also observe that in ResNet-18, the number of communication rounds grows faster than expected for  
 319 larger batch sizes. The reason for this is that the Ripple-carry adder used in the conversion from  $[x]$   
 320 to  $\langle x \rangle$  is very memory-intensive. When CRYPTEN runs out of GPU memory, it replaces the adder  
 321 by an implementation that requires  $O(|\mathcal{P}|)$  communication rounds (compared to  $(\log_2 |\mathcal{P}|)$  for the  
 322 Ripple-carry adder) but that requires less memory (see appendix).

<sup>6</sup>We adopted the ResNet implementation from torchvision and the ViT implementation from <https://github.com/rwightman/pytorch-image-models>.

## 323 7 Conclusion and Future Work

324 In this paper, we have introduced and benchmarked CRYPTEN. We hope that CRYPTEN’s flexible,  
325 machine-learning first API design and performance can help foster adoption of secure MPC in  
326 machine learning. We see the following directions for future research and development of CRYPTEN.

327 **Numerical issues** are substantially more common in CRYPTEN implementations of machine-learning  
328 algorithms than in their PyTorch counterparts. In particular, the fixed-point representation with  $L$   
329 of precision ( $L=16$  by default) is more prone to numerical overflow or underflow than floating-point  
330 representations. Moreover, arithmetic secret shares are prone to *wrap-around* errors in which the sum  
331 of the shares  $[x]_p$  exceeds the size of the ring,  $Q=2^{64}$ . Wrap-around errors can be difficult to debug  
332 because they may only arise in the multi-party setting, in which no individual party can detect them.  
333 We plan to implement tools in CRYPTEN that assist users in debugging such numerical issues.

334 **End-to-end privacy** requires seamless integration between data-processing frameworks, such as  
335 secure SQL implementations [5], and data-modeling frameworks like CRYPTEN. In “plaintext”  
336 software, such frameworks are developed independently and combined via “glue code” or platforms  
337 that facilitate the construction of processing and modeling pipelines. Real-world use cases of machine  
338 learning via secure MPC require the development of a platform that makes the integration of private  
339 data processing and modeling seamless, both from an implementation and a security point-of-view.

340 **Differential privacy** mechanisms may be required in real-world applications of CRYPTEN in order  
341 to provide rigorous guarantees on the information leakage that inevitably occurs when the results of a  
342 private computation are publicly revealed [27]. CRYPTEN implements sampling algorithms for the  
343 Bernoulli, Laplace, and Gaussian distributions (see appendix), which allows for the implementation  
344 of randomized response [67], the Laplace mechanism [28], and the Gaussian mechanism [6, 27]  
345 (although care must be taken when implementing these mechanisms [12, 45]). In future work, we  
346 aim to use these mechanisms, for example, to do a secure MPC implementation of DP-SGD [2].

347 **Threat models** may vary per use case. Specifically, some use cases may require malicious security or  
348 may not provide a TTP. We plan to add support for malicious security via message authentication codes  
349 to CRYPTEN [22]. We also plan to support generation of Beaver triples via additive homomorphic  
350 encryption [51], oblivious transfer [38], or more recent methods [10] to eliminate the need for a TTP.

351 **Model architecture design** for secure MPC is another important direction for future research.  
352 Following prior work in this research area, this study has focused on implementing *existing* machine-  
353 learning models in a secure MPC framework. However, these models were designed based on  
354 computational considerations in “plaintext” implementations of the models on modern GPU or  
355 TPU hardware. The results of our benchmarks suggest that this may be suboptimal because those  
356 considerations are very different in a secure MPC environment. For example, the evaluation of softmax  
357 functions over large numbers of values requires a lot of communication in secure MPC, which makes  
358 attention layers very slow. This implies that multilayer perceptron models [63] are likely much  
359 more efficient than vision transformers [25, 64] for image classification. We hope that CRYPTEN’s  
360 machine-learning API and ease of use will spur studies that design model architectures specifically  
361 optimized for a secure MPC environment, for example, via neural architecture search [43, 46, 70].

## 362 8 Broader Impact

363 Although we believe that the adoption of secure MPC in machine learning can lead to the development  
364 of AI systems that are substantially more private and secure, we note that there are also potential  
365 downsides to such adoption. In particular, because the computations in secure MPC are performed on  
366 encrypted data, it can be harder to do quality control of AI systems implemented in CRYPTEN. For  
367 example, it is impossible to inspect the values of intermediate activations (or even model outputs)  
368 unless all parties agree to reveal those values. This may make it harder to explain why a model  
369 makes a certain decision [24] or to detect data-poisoning attacks [8]. Indeed, there exist fundamental  
370 trade-offs between privacy and utility [56] and those trade-offs apply to CRYPTEN users, too.

371 It is also worth noting that, although the protocols implemented in CRYPTEN come with rigorous  
372 cryptographic guarantees, practical implementations of these protocols may be broken by other means.  
373 For example, we have no reason to assume that CRYPTEN would not be susceptible to side-channel  
374 attacks [60]. Hence, good data stewardship remains essential even when using secure computation.

375 **References**

- 376 [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis,  
377 J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Joze-  
378 fowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah,  
379 M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Va-  
380 sudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng.  
381 TensorFlow: Large-scale machine learning on heterogeneous systems, 2015.
- 382 [2] M. Abadi, A. Chu, I. Goodfellow, H. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep  
383 learning with differential privacy. In *Proceedings of the CCS*, pages 308–318, 2016.
- 384 [3] A. Agrawal, A. N. Modi, A. Passos, A. Lavoie, A. Agarwal, A. Shankar, I. Ganichev, J. Leven-  
385 berg, M. Hong, R. Monga, and S. Cai. TensorFlow Eager: A multi-stage, python-embedded dsl  
386 for machine learning. In *arXiv:1903.01855*, 2019.
- 387 [4] T. Araki, J. Furukawa, Y. Lindell, A. Nof, and K. Ohara. High throughput semi-honest secure  
388 three-party computation with an honest majority. In *ACM CCS*, pages 805–817, 2016.
- 389 [5] D. W. Archer, D. Bogdanov, L. Kamm, Y. Lindell, K. Nielsen, J. I. Pagter, N. P. Smart, and R. N.  
390 Wright. From keys to databases – real-world applications of secure multi-party computation.  
391 *Computer Journal*, 61(12):1749–1771, 2018.
- 392 [6] B. Balle and Y.-X. Wang. Improving the Gaussian mechanism for differential privacy. In  
393 *Proceedings of International Conference on Machine Learning*, 2018.
- 394 [7] D. Beaver. Efficient multiparty protocols using circuit randomization. In *Annual International*  
395 *Cryptology Conference*, pages 420–432. Springer, 1991.
- 396 [8] B. Biggio, B. Nelson, and P. Laskov. Poisoning attacks against support vector machines. In  
397 *International Conference on Machine Learning (ICML)*, pages 1467–1474, 2012.
- 398 [9] D. Bogdanov, S. Laur, and J. Willemson. Sharemind: A framework for fast privacy-preserving  
399 computations. In S. Jajodia and J. Lopez, editors, *Computer Security - ESORICS 2008*, pages  
400 192–206. Springer Berlin Heidelberg, 2008.
- 401 [10] E. Boyle, G. Couteau, N. Gilboa, Y. Ishai, L. Kohl, and P. Scholl. Correlated pseudorandom  
402 functions from variable-density LPN. In *Cryptology ePrint Archive: Report 2020/1417*, 2020.
- 403 [11] R. Canetti. Universally composable security: A new paradigm for cryptographic protocols. In  
404 *Proceedings of the Annual Symposium on Foundations of Computer Science*, pages 136–145,  
405 2001.
- 406 [12] C. Canonne, G. Kamath, and T. Steinke. The discrete Gaussian for differential privacy. In *arXiv*  
407 *2004.00010*, 2020.
- 408 [13] Cape Privacy. TF-Trusted. URL <https://github.com/capeprivacy/tf-trusted>.
- 409 [14] O. Catrina and S. De Hoogh. Improved primitives for secure multiparty integer computation. In  
410 *International Conference on Security and Cryptography for Networks*, pages 182–199. Springer,  
411 2010.
- 412 [15] N. Chandran, D. Gupta, A. Rastogi, R. Sharma, and S. Tripathi. EzPC: Programmable, efficient,  
413 and scalable secure two-party computation for machine learning. In *IEEE European Symposium*  
414 *on Security and Privacy*, 2019.
- 415 [16] H. Chaudhari, A. Choudhury, A. Patra, and A. Suresh. Astra: High throughput 3pc over  
416 rings with application to secure prediction. In *ACM SIGSAC Conference on Cloud Computing*  
417 *Security Workshop*, 2019.
- 418 [17] R. Collobert, C. Puhersch, and G. Synnaeve. Wav2letter: An end-to-end convnet-based speech  
419 recognition system. In *arXiv:1609.03193*, 2016.

- 420 [18] R. Cramer, I. Damgård, and Y. Ishai. Share conversion, pseudorandom secret-sharing and  
421 applications to secure computation. In *Lecture Notes in Computer Science*, volume 3378, pages  
422 342–362, 2005.
- 423 [19] cuBLAS. <https://developer.nvidia.com/cublas>.
- 424 [20] cuDNN. <https://developer.nvidia.com/cudnn>.
- 425 [21] I. Damgård, M. Fitzi, E. Kiltz, J. B. Nielsen, and T. Toft. Unconditionally secure constant-rounds  
426 multi-party computation for equality, comparison, bits and exponentiation. In *TCC*, 2005.
- 427 [22] I. Damgård, V. Pastro, N. Smart, and S. Zakarias. Multiparty computation from somewhat  
428 homomorphic encryption. Cryptology ePrint Archive, Report 2011/535, 2011. [https://](https://eprint.iacr.org/2011/535)  
429 [eprint.iacr.org/2011/535](https://eprint.iacr.org/2011/535).
- 430 [23] D. Demmler, T. Schneider, and M. Zohner. ABY – A framework for efficient mixed-protocol  
431 secure two-party computation. In *NDSS*, 2015.
- 432 [24] F. Doshi-Velez and B. Kim. Towards a rigorous science of interpretable machine learning. In  
433 *arXiv:1702.08608*, 2017.
- 434 [25] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,  
435 M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16  
436 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- 437 [26] T. Dugan and X. Zou. A survey of secure multiparty computation protocols for privacy  
438 preserving genetic tests. In *Proceedings of the International Symposium on Biomedical Imaging*,  
439 2016.
- 440 [27] C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Foundations and*  
441 *Trends in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- 442 [28] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private  
443 data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- 444 [29] D. Evans, V. Kolesnikov, and M. Rosulek. *A Pragmatic Introduction to Secure Multi-Party*  
445 *Computation*. NOW Publishers, 2018.
- 446 [30] Gloo. <https://github.com/facebookincubator/gloo>.
- 447 [31] O. Goldreich, S. Micali, and A. Wigderson. How to play any mental game or a completeness  
448 theorem for protocols with honest majority. In *STOC*, pages 218–229, 1987.
- 449 [32] V. Haralampieva, D. Rueckert, and J. Passerat-Palmbach. A systematic comparison of encrypted  
450 machine learning solutions for image classification. In *Proceedings of the Workshop on Privacy-*  
451 *Preserving Machine Learning in Practice*, pages 55–59, 2020.
- 452 [33] M. Hastings, B. Hemenway, D. Noble, and S. Zdancewic. SoK: general-purpose compilers for  
453 secure multi-party computation. In *2019 IEEE Symposium on Security and Privacy (SP)*, 2019.
- 454 [34] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In  
455 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,  
456 2016.
- 457 [35] A. S. Householder. *The Numerical Treatment of a Single Nonlinear Equation*. 1970.
- 458 [36] M. Ito, A. Saito, and T. Nishizeki. Secret sharing scheme realizing general access structure.  
459 *Electronics and Communications in Japan (Part III: Fundamental Electronic Science)*, 72:  
460 56–64, 1989.
- 461 [37] C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan. Gazelle: A low latency framework for  
462 secure neural network inference. In *arXiv 1801.05507*, 2018.
- 463 [38] M. Keller, E. Orsini, and P. Scholl. MASCOT: Faster malicious arithmetic secure computation  
464 with oblivious transfer. In *Proceedings of the ACM SIGSAC Conference on Computer and*  
465 *Communications Security*, pages 830–842, 2016.

- 466 [39] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional  
467 neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- 468 [40] N. Kumar, M. Rathee, N. Chandran, D. Gupta, A. Rastogi, and R. Sharma. CryptFlow: Secure  
469 TensorFlow inference. In *IEEE Symposium on Security and Privacy*, pages 336–353, 2020.
- 470 [41] A. Lapets, N. Volgushev, A. Bestavros, F. Jansen, and M. Varia. Secure multi-party computation  
471 for analytics deployed as a lightweight web application. Technical Report BU-CS-TR 2016-008,  
472 Boston University, 2016.
- 473 [42] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document  
474 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 475 [43] H. Liu, K. Simonyan, and Y. Yang. DARTS: Differentiable architecture search. In  
476 *arXiv:1806.09055*, 2018.
- 477 [44] J. Liu, M. J. Y. Lu, and N. Asokan. Oblivious neural network predictions via MiniONN  
478 transformations. In *ACM Conference on Computer and Communications Security (CCS)*, 2017.
- 479 [45] I. Mironov. On significance of the least significant bits for differential privacy. In *Proceedings*  
480 *ACM Conference on Computer and Communications Security (CCS)*, pages 650–661, 2012.
- 481 [46] P. Mishra, R. Lehmkuhl, A. Srinivasan, W. Zheng, and R. Popa. Delphi: A cryptographic  
482 inference service for neural networks. In *USENIX Security Symposium*, 2020.
- 483 [47] P. Mohassel and P. Rindal. ABY3: A mixed protocol framework for machine learning. In *ACM*  
484 *Conference on Computer and Communications Security (CCS)*, 2018.
- 485 [48] P. Mohassel and Y. Zhang. SecureML: A system for scalable privacy-preserving machine  
486 learning. In *IEEE Symposium on Security and Privacy*, 2017.
- 487 [49] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In  
488 *Proceedings of the International Conference on Machine Learning (ICML)*, 2010.
- 489 [50] NCCL. <https://developer.nvidia.com/nccl>.
- 490 [51] P. Paillier. Public-key cryptosystems based on composite degree residuosity classes. In  
491 *EUROCRYPT*, pages 223–238, 1999.
- 492 [52] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur. LibriSpeech: An ASR corpus based  
493 on public domain audio books. In *Proceedings of the International Conference on Acoustics,*  
494 *Speech, and Signal Processing (ICASSP)*, 2015.
- 495 [53] A. Paszke, S. Gross, S. Chintala, and G. Chanan. PyTorch: Tensors and dynamic neural  
496 networks in Python with strong GPU acceleration, 2017.
- 497 [54] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,  
498 N. Gimeshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani,  
499 S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. PyTorch: An imperative style,  
500 high-performance deep learning library. In *Advances in Neural Information Processing Systems*  
501 *(NeurIPS)*, 2019.
- 502 [55] A. Patra and A. Suresh. Blaze: Blazing fast privacy preserving machine learning. In *Symposium*  
503 *on Network and Distributed System Security (NDSS)*, 2020.
- 504 [56] I. S. Reed. Information theory and privacy in data banks. In *Proceedings of the June 4-8, 1973,*  
505 *National Computer Conference and Exposition*, 1973.
- 506 [57] M. S. Riazi, C. Weinert, O. Tkachenko, E. M. Songhori, T. Schneider, and F. Koushanfar.  
507 Chameleon: A hybrid secure computation framework for machine learning applications. In  
508 *Cryptology ePrint Archive*, volume 2017/1164, 2017.
- 509 [58] M. S. Riazi, M. Samragh, H. Chen, K. Laine, K. Lauter, and F. Koushanfar. Xonn: Xnor-based  
510 oblivious deep neural network inference. In *USENIX Security*, 2019.

- 511 [59] T. Ryffel, A. Trask, M. Dahl, B. Wagner, J. Mancuso, D. Rueckert, and J. Passerat-Palmbach. A  
512 generic framework for privacy preserving deep learning. In *arXiv:1811.04017*, 2018.
- 513 [60] O. Seker, S. Berndt, L. Wilke, and T. Eisenbarth. SNI-in-the-head: Protecting MPC-in-the-head  
514 protocols against side-channel analysis. In *Cryptology ePrint Archive: Report 2020/544*, 2020.
- 515 [61] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image  
516 recognition. In *ICLR*, 2015.
- 517 [62] S. Tan, B. Knott, Y. Tian, and D. J. Wu. CryptGPU: Fast privacy-preserving machine learning  
518 on the GPU. In *arXiv 2104.10949*, 2021.
- 519 [63] I. Tolstikhin, N. Houlsby, A. Kolesnikov, L. Beyer, X. Zhai, T. Unterthiner, J. Yung, A. Steiner,  
520 D. Keysers, J. Uszkoreit, M. Lucic, and A. Dosovitskiy. MLP-Mixer: An all-MLP architecture  
521 for vision. In *arXiv:2105.01601*, 2021.
- 522 [64] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and  
523 I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*  
524 (*NeurIPS*), 2017.
- 525 [65] S. Wagh, D. Gupta, and N. Chandran. SecureNN: Efficient and private neural network training.  
526 In *Cryptology ePrint Archive*, volume 2018/442, 2018.
- 527 [66] S. Wagh, S. Tople, F. Benhamouda, E. Kushilevitz, P. Mittal, and T. Rabin. FALCON: Honest-  
528 majority maliciously secure framework for private deep learning. In *Proc. Priv. Enhancing*  
529 *Technol.*, 2021.
- 530 [67] S. L. Warner. Randomized response: A survey technique for eliminating evasive answer bias.  
531 *Journal of the American Statistical Association*, 60(309):63–69, 1965.
- 532 [68] A. C.-C. Yao. How to generate and exchange secrets. In *FOCS*, pages 162–167, 1986.
- 533 [69] Yelp. Yelp Review Dataset. URL <https://www.yelp.com/dataset>.
- 534 [70] B. Zoph and Q. V. Le. Neural architecture search with reinforcement learning. In  
535 *arXiv:1611.01578*, 2016.

536 **Checklist**

- 537 1. For all authors...
- 538 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
539 contributions and scope? [Yes]
- 540 (b) Did you describe the limitations of your work? [Yes] See Section 7.
- 541 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See  
542 Section 8.
- 543 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
544 them? [Yes]
- 545 2. If you are including theoretical results...
- 546 (a) Did you state the full set of assumptions of all theoretical results? [N/A] We do not  
547 present new theoretical results.
- 548 (b) Did you include complete proofs of all theoretical results? [N/A] The security proof  
549 follows trivially from prior work. The appendix contains pointers to the relevant proofs.
- 550 3. If you ran experiments...
- 551 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
552 mental results (either in the supplemental material or as a URL)? [Yes] See Section 6.
- 553 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
554 were chosen)? [N/A] We do not train any models in this paper.
- 555 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
556 ments multiple times)? [No] The standard deviations in the results of our experiments  
557 are negligible or zero, which is why we do not include error bars.
- 558 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
559 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 6.
- 560 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 561 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 6.
- 562 (b) Did you mention the license of the assets? [No] All code we used is licensed under  
563 Apache or BSD licenses. CRYPTEN itself is licensed under an MIT license. The  
564 LibriSpeech dataset has a Creative Commons (CC-BY 4.0) license. The license of the  
565 ImageNet dataset is unclear. The Yelp Review dataset has a custom license.
- 566 (c) Did you include any new assets either in the supplemental material or as a URL?  
567 [No] The CRYPTEN library is available as open-source, but we did not include it as  
568 supplemental material or URL because doing so would reveal information about the  
569 identity of the authors.
- 570 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
571 using/curating? [N/A] We do not use data on people in this study.
- 572 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
573 information or offensive content? [No] The Yelp Review, LibriSpeech, or ImageNet  
574 dataset may contain personally identifiable information or offensive content, but the  
575 data content is not the subject of this study. Specifically, we have not used these datasets  
576 to train models.
- 577 5. If you used crowdsourcing or conducted research with human subjects...
- 578 (a) Did you include the full text of instructions given to participants and screenshots, if  
579 applicable? [N/A] We did not conduct research with human subjects.
- 580 (b) Did you describe any potential participant risks, with links to Institutional Review  
581 Board (IRB) approvals, if applicable? [N/A] We did not conduct research with human  
582 subjects.
- 583 (c) Did you include the estimated hourly wage paid to participants and the total amount  
584 spent on participant compensation? [N/A] We did not conduct research with human  
585 subjects.