

# Modular Sumcheck Proofs with Applications to Machine Learning and Image Processing

---

**D. Balbás**<sup>1,2</sup>, D. Fiore<sup>1</sup>, M. González-Vasco<sup>3</sup>, D. Robissout<sup>1</sup>, C. Soriente<sup>4</sup>

28th November 2023

<sup>1</sup>IMDEA Software Institute, Madrid, Spain

<sup>2</sup>Universidad Politécnica de Madrid, Spain

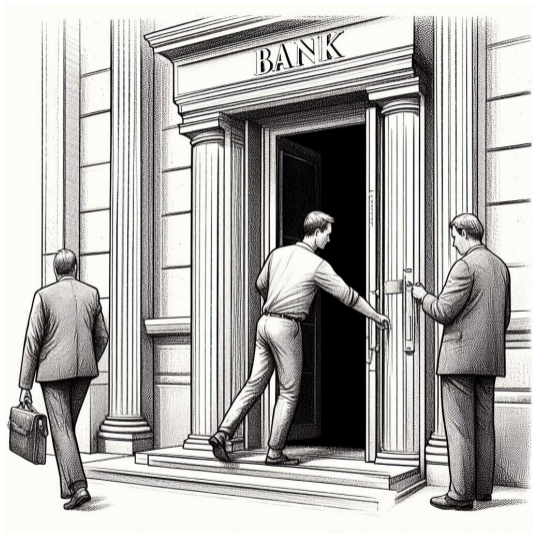
<sup>3</sup>Universidad Carlos III de Madrid, Spain

<sup>4</sup>NEC Laboratories Europe, Madrid, Spain

ACM CCS 2023, Copenhagen, Denmark







# Proof Systems in the Wild

Towards proving larger models, we require:

- **Efficient Verification**
- **Efficient Proof Generation:**  $\tilde{O}(n)$  time, usually achieved by *sumcheck-based proofs*. Low memory usage.
- **Privacy** for model parameters.

# Proof Systems in the Wild

Towards proving larger models, we require:

- **Efficient Verification**
- **Efficient Proof Generation:**  $\tilde{O}(n)$  time, usually achieved by *sumcheck-based proofs*. Low memory usage.
- **Privacy** for model parameters.

**General-purpose** proof systems (e.g. SNARKs) not great for *data intensive computations*.

# Proof Systems in the Wild

Towards proving larger models, we require:

- **Efficient Verification**
- **Efficient Proof Generation:**  $\tilde{O}(n)$  time, usually achieved by *sumcheck-based proofs*. Low memory usage.
- **Privacy** for model parameters.

**General-purpose** proof systems (e.g. SNARKs) not great for *data intensive computations*.

**Special-purpose** proofs (e.g. vCNN, zkCNN, zkIMG) better but lack *composability/reusability*.

# This Work

- **Framework** for composing sumcheck-based proofs at an information-theoretic level.

# This Work

- **Framework** for composing sumcheck-based proofs at an information-theoretic level.
- Better efficiency for **combined special-purpose protocols**.



# This Work

- **Framework** for composing sumcheck-based proofs at an information-theoretic level.
- Better efficiency for **combined special-purpose protocols**.
- Efficient proofs for multi-channel **convolution**.

# This Work

- **Framework** for composing sumcheck-based proofs at an information-theoretic level.
- Better efficiency for **combined special-purpose protocols**.
- Efficient proofs for multi-channel **convolution**.
- Modular, extendable Rust **implementation**, 5-10x faster & shorter than special-purpose proofs for ML and IP.

This work as seen by Dall·E 3

# This Work

- **Framework** for *composing sumcheck-based proofs* at an information-theoretic level.
- Better efficiency for **combined special-purpose protocols**.
- Efficient *proofs for multi-channel convolution*.
- Modular, extendable Rust **implementation**, 5-10x *faster & shorter than special-purpose proofs* for ML and IP.

This work as seen by Dall·E 3

## Our Framework

---

# Fingerprints and Interactive Proofs

An **interactive proof** (IP) for the language  $L_F = \{f(x) : f(x) = yg\}$  is a pair of algorithms  $(P; V)$  that are *complete* and *sound*.

# Fingerprints and Interactive Proofs

An **interactive proof** (IP) for the language  $L_F = \{f(x) : f(x) = yg\}$  is a pair of algorithms  $(P, V)$ ;  $\forall (f; x; y) \in L_F$  that are *complete* and *sound*.

## Fingerprint

Let  $c_x = H(x; r)$  be the **fingerprint** of  $x$  on  $r$ .  
 $H$  is (statistically) sound if for *any* pair  $x \neq x'$ ,

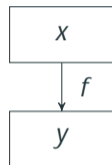
$$\Pr_r[H(x; r) = H(x'; r)] = \text{negl}(\lambda):$$

Example: for  $x \in F^n$ , the poly. evaluation  $H(x; r) = x_0 + x_1r + \dots + x_{n-1}r^{n-1}$  over  $F$ .

# Structure of Common IPs

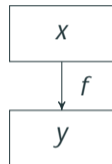
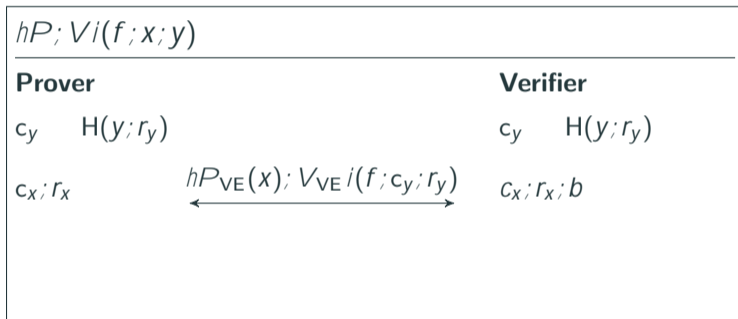
We can explain many efficient IPs for  $(f; x; y)$  with the following abstraction:

$hP; Vi(f; x; y)$	
<b>Prover</b>	<b>Verifier</b>
$c_y \quad H(y; r_y)$	$c_y \quad H(y; r_y)$



# Structure of Common IPs

We can explain many efficient IPs for  $(f; x; y)$  with the following abstraction:



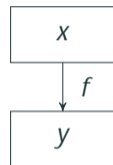


# Structure of Common IPs

We can explain many efficient IPs for  $(f; x; y)$  with the following abstraction:

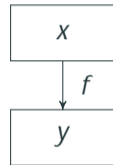
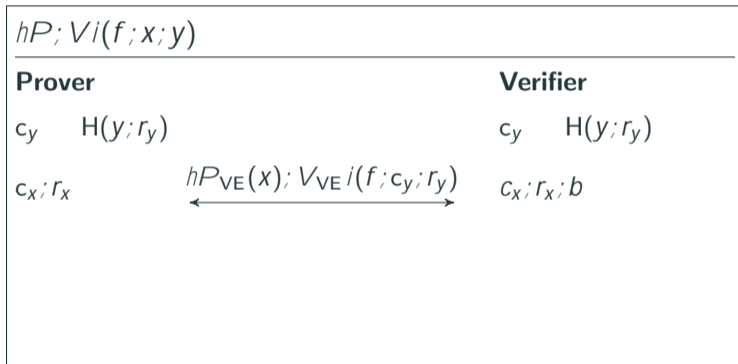
$hP; Vi(f; x; y)$	
<b>Prover</b>	<b>Verifier</b>
$c_y \quad H(y; r_y)$	$c_y \quad H(y; r_y)$
$c_x; r_x$	$c_x; r_x; b$
	$b^\theta \quad [c_x = H(x; r_x)]$
	<b>return</b> $b \wedge b^\theta$

$\xleftrightarrow{hP_{VE}(x); V_{VE}i(f; c_y; r_y)}$



# Structure of Common IPs

We can explain many efficient IPs for  $(f; x; y)$  with the following abstraction:

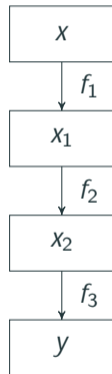


Subroutines named **verifiable evaluation schemes (VE)** on *ngerprinted data*.

# Structure of Common IPs

We can explain many efficient IPs for  $(f; x; y)$  with the following abstraction:

$hP; Vi(f; x; y)$		
Prover		Verifier
$c_y$	$H(y; r_y)$	$c_y$
$c_2; r_2$	$\xleftrightarrow{hP_{VE}(x_2); V_{VE}i(f_3; c_y; r_y)}$	$c_2; r_2; b_2$
$c_1; r_1$	$\xleftrightarrow{hP_{VE}(x_1); V_{VE}i(f_2; c_2; r_2)}$	$c_1; r_1; b_1$
$c_x; r_x$	$\xleftrightarrow{hP_{VE}(x); V_{VE}i(f_1; c_1; r_1)}$	$c_x; r_x; b_0$



Subroutines named **verifiable evaluation schemes (VE)** on *ngerprinted data*.

# Our Framework

We characterize VEs and provide a formalism. **VEs can be composed sequentially at the information-theoretic level!**

VEs can express many sumcheck-based proof systems:

# Our Framework

We characterize VEs and provide a formalism. **VEs can be composed sequentially at the information-theoretic level!**

VEs can express many sumcheck-based proof systems:

- Matrix multiplication [Thaler13]
- GKR [GKR08, CMT12]
- Libra [XZZPS19], Virgo [ZLW+20] (and follow-ups)
- FFT-based convolution [LXZ21]
- ...

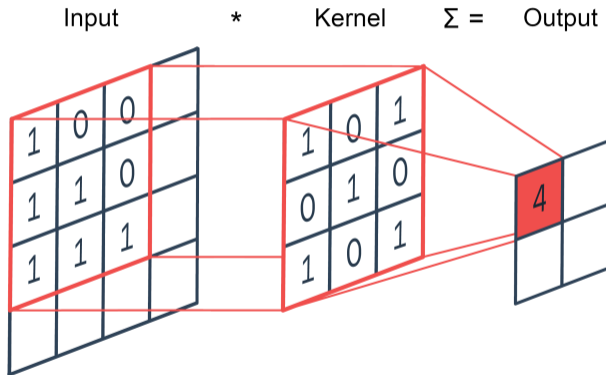
We provide a compiler *from VEs to succinct cryptographic arguments*. Fingerprints can be evaluated efficiently by  $\mathcal{V}$  via polynomial commitments.

# Efficient Proofs for Convolution

---

# Our Proofs for Convolution

We express convolutions between input  $X$  and kernel  $W$  as multilinear sumchecks.



Source: Christopher Melen, RNCM

# Our Proofs for Convolution

Our protocol proceeds in two steps:

1. A *reshape* sumcheck that rearranges  $X \nabla \hat{X}$ .
2. A *convolution* sumcheck  $\hat{X} \cdot W \nabla Y$ .

Built upon sumchecks for matrix multiplication [Tha13] and channel batching.



# Our Proofs for Convolution

Our protocol proceeds in two steps:

1. A *reshape* sumcheck that rearranges  $X \not\approx \hat{X}$ .
2. A *convolution* sumcheck  $\hat{X} \cdot W \not\approx Y$ .

Built upon sumchecks for matrix multiplication [Tha13] and channel batching.

## Performance

For  $c$  input channels,  $d$  output channels,

	Ours	zkCNN [LXZ21]
<b>Prover</b>	$O(c \cdot  W  \cdot ( Y  + d))$	$O(c \cdot d \cdot  X )$
<b>Verifier</b>	$O(\log(c \cdot  Y ))$	$O(\log^2(c \cdot d \cdot  X ))$
<b>Size</b>	$O(\log(c \cdot  Y ))$	$O(\log^2(c \cdot d \cdot  X ))$

# Applications and Benchmarking

---

We extend our framework to construct efficient proof systems for:

- **Convolutional Neural Networks.**
- **Recurrent NNs.**
- **Image Processing:** Native linear, reshaping, and convolutional operations (filtering, blurring...).

Our convolution prover and a general CNN prover are implemented in Rust and available open-source.

# Benchmarking

- **Prover** 0.1s for 256 × 256 input and 4 × 4 kernel.  
5 × faster than zkCNN, 100 × faster than vCNN.

# Benchmarking

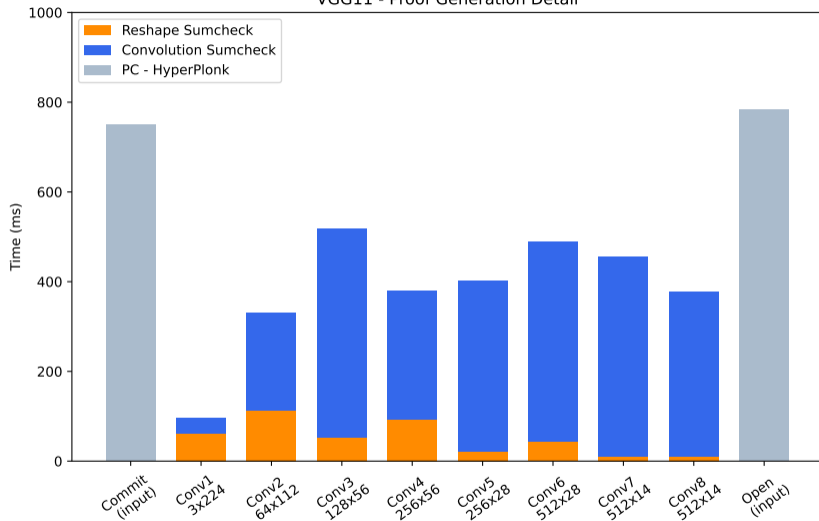
- **Prover** 0:1s for 256 × 256 input and 4 × 4 kernel.  
5 *faster than zkCNN*, 100 *faster than vCNN*.
- **Verification** 0:1ms.

# Benchmarking

- **Prover** 0:1s for 256 × 256 input and 4 × 4 kernel.  
5 *faster than zkCNN*, 100 *faster than vCNN*.
- **Verification** 0:1ms.
- **Proof size** 1KB.  
10 *shorter than zkCNN*.

*Sequential composition reduces memory usage.* Still room for improvement.

VGG11 - Proof Generation Detail



All kernels are  $3 \times 3$ . Run on single-core Xeon-Gold-6154 at 3GHz.

### VGG11

- Conv1: 3  $224^2$
- Conv2: 64  $112^2$
- Conv3: 128  $56^2$
- Conv4: 256  $56^2$
- Conv5: 256  $28^2$
- Conv6: 512  $28^2$
- Conv7: 512  $14^2$
- Conv8: 512  $14^2$

# Conclusions

- Theory **framework** for composition of sumcheck-based proofs.
- Efficient sumchecks for **convolution**.
- Efficient arguments for **data-intensive applications** such as ML and IP.
- Modular **implementation** in Rust.



# Conclusions

- Theory **framework** for composition of sumcheck-based proofs.
- Efficient sumchecks for **convolution**.
- Efficient arguments for **data-intensive applications** such as ML and IP.
- Modular **implementation** in Rust.

*Thank you!*

`ia.cr/2023/1342`

`david.balbas@imdea.org`

*Slides available @ `davidbalbas.github.io`*

## **Additional Material**

---

# VE for Multilinear Polynomials

Let  $t = (t_1; \dots; t_n)$ ,  $F$  a finite field, and

$$x(t; y) = \prod_{i=1}^s x_i(t; y)$$

where each  $x_i$  is a multilinear polynomial over  $F$ .

We can generalize multilinear sumcheck as a VE for the relation

$$f_y(r_y) = \prod_{t \in \mathbb{F}^n} x(t; r_y):$$

$P; V$  start on fingerprint  $c_y = f_y(r_y)$ . At the end, they obtain  $c_x = x(r_t; r_y)$ .

# VE for Matrix Multiplication [Tha13]

Let  $C = A \cdot B$  where  $A, B, C \in \mathbb{F}^{n \times n}$ . We can write matrix multiplication as

$$\tilde{C}(x_1; x_2) = \sum_{y \in \mathbb{F}^n} \tilde{A}(x_1; y) \tilde{B}(y; x_2)$$

Where  $\tilde{A}$  encodes  $A$  as a (unique) multilinear polynomial,  $\tilde{A}(i; j) = A_{ij}$

Given  $r_1, r_2 \in \mathbb{F}^n$ , we apply our multilinear sumcheck VE on:

$$\tilde{C}(r_1; r_2) = \sum_{y \in \mathbb{F}^n} \tilde{A}(r_1; y) \tilde{B}(y; r_2)$$

$P; V$  start on *ngerprint*  $c_C = \tilde{C}(r_1; r_2)$ . At the end, they *obtain ngerprints*  $c_A = \tilde{A}(r_1; r_3)$  and  $c_B = \tilde{B}(r_3; r_2)$ .

Communication and verifier  $t_V = O(\log n) = O(\log n)$ , prover  $O(n^2)$ .

# VE for Convolution

Convolution equations can be compacted as

$$\text{vec}(Y) = \begin{matrix} & \begin{matrix} 2 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} W_0 X_0 + W_1 X_1 + W_3 X_3 + W_4 X_4 \\ W_0 X_1 + W_1 X_2 + W_3 X_4 + W_4 X_5 \\ W_0 X_3 + W_1 X_4 + W_3 X_6 + W_4 X_7 \\ W_0 X_4 + W_1 X_5 + W_3 X_7 + W_4 X_8 \end{matrix} & \begin{matrix} 3 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} X_0 \\ X_1 \\ X_3 \\ X_4 \\ X_5 \\ X_4 \end{matrix} & \begin{matrix} 2 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} X_1 \\ X_2 \\ X_4 \\ X_4 \\ X_5 \end{matrix} & \begin{matrix} 3 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} X_3 \\ X_4 \\ X_6 \\ X_7 \\ X_7 \end{matrix} & \begin{matrix} 3 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} X_4 \\ X_5 \\ X_6 \\ X_7 \\ X_8 \end{matrix} & \begin{matrix} 3 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} W_0 \\ W_1 \\ W_3 \\ W_3 \\ W_4 \end{matrix} \end{matrix}$$

For **multiple kernels and inputs** (multi-channel), as usual in e.g. neural networks,

$$Y = [Y_{1j} \quad jY_d] = \prod_{=1}^c \hat{X} [W_{,1j} \quad jW_{,d}]$$

Where  $2[c]$  represents input channels, and  $2[d]$  output channels.