Hyperplonk: Plonk with linear-time prover and high degree gates

Binyi Chen Espresso Systems Benedikt Bünz Espresso, Stanford, NYU Dan Boneh Stanford Zhenfei Zhang Espresso Systems

Preprocessing (zk)-SNARKs

SNARK = A succinct proof showing that $\exists w$ s.t. C(x, w) = 0small size fast to verify

Why so popular recently? Outsourcing computation





Computation happens locally



Computation goes remote/global



Application: Blockchain

SNARK = A succinct proof showing that $\exists w$ s.t. C(x, w) = 0



More applications: zkRollup, zkEVM, zkBridge, DSNs, zkML

Preprocessing (zk)-SNARKs

Call for:

- Fast prover for large statements (e.g., as fast as computation)
- Small proof size and efficient verification
- Powerful circuit constraint systems (e.g., high-deg/lookup gates)
- Hardware-friendliness

Our Contributions

HyperPlonk+ PolyIOP

- Linear-time prover, no use of FFT-friendly field
 - The *first* linear-time SNARK that has expressive gate support and small proof size
 - Hardware-optimization friendly
- Native high-degree custom gate support
 - Better support than Plonk [GWC19]
- Lookup gate support
 - The *first linear-time* PolyIOP for the lookup relation!
- Towards strict linear-time SNARKs: Orion+ PCS
 - Improve opening size of the state-of-the-arts (Orion & Brakedown)
 - 5MB proof -> <10KB proof
- Active industry deployment/development
 - E.g., ETH Foundation/Scroll actively develop it for future zkEVM solutions



General paradigm for modern SNARKs

Many SNARKs are built in two steps:









Binding : cannot output two valid openings (f_1, r_1) , (f_2, r_2) for com _f .	
Optional: commitment com_f is hiding	

Polynomial Commitment



Optional: **Eval** is **zero knowledge**: π "reveals nothing" about f.

Outline

- A generic framework for proving circuit relations
- High-degree gates support
- Hyperplonk+: support lookup on the Boolean hypercube
- Evaluations

Outline

- A generic framework for proving circuit relations
- High-degree gates support
- Hyperplonk+: support lookup on the Boolean hypercube
- Evaluations

Workflow [GWC19]



PLONK's Circuit Encoding



Trace to Polynomial

The computation trace:

Gate 0	Gate 1	Gate 2	Gate 3
5	6	11	77
6	1	7	7
11	7	77	84

The polynomials:

L(H) =	5	6	11	77
R(H) =	6	1	7	7
O(H) =	11	7	77	84

Options of H:

- $\mathbf{H} = \{1, \omega, \omega^2, \omega^3\}$
 - L, R, O are univariate polys
 - Interpolating L, R, O needs FFT
- **H** = {00,01,10,11}
 - L, R, O are multilinear polys
 - Polys are in eval form, no interpolation
 - Free embedding!

Gate Check [GWC19]

The polynomials:

L(H) =	5	6	11	77
R(H) =	6	1	7	7
O(H) =	11	7	77	84
S(H) =	1	1	0	1



Idea: encode gate types using a <u>selector</u> polynomial S(X)

$$\forall \ell = 0, ..., |C| - 1$$
:
 $S(H_{\ell}) = 1$ if gate #*l* is an addition gate
 $S(H_{\ell}) = 0$ if gate #*l* is a multiplication gate

Gate Check [GWC19]

The polynomials:

L(<mark>H</mark>) =	5	6	11	77
R(H) =	6	1	7	7
O(H) =	11	7	77	84
S(H) =	1	1	0	1



The zero-check: $\forall y \in H$:

$$S(y) \cdot [L(y) + R(y)] + (1 - S(y)) \cdot [L(y) \cdot R(y)] = O(y)$$

Addition gate

Multiplication gate

- $1 \cdot [L(y) + R(y)] + (1 1) \cdot = O(y) \qquad 0 \cdot + (1 0) \cdot [L(y) \cdot R(y)] = O(y)$
- 0(y)

Gate Check [GWC19]

The polynomials:

L(H) =	5	6	11	77
R(H) =	6	1	7	7
O(H) =	11	7	77	84
S(H) =	1	1	0	1



SumCheck $\sum_{y \in H} f(y) \cdot eq_r(y) = 0$

Intuition: rand linear combine

Univariate Sumcheck: requires FFTs

ZeroCheck

 $f(y) = 0 \forall y \in H$

Multivariate Sumcheck: linear prover time!

Wiring Identity Check



The computation trace:



All slots with the same color should have identical values!

(L, R, O) $\pi(L, R, O)$

for some **fixed permutation** π

Permutation Check

Prove $(L, R, O) = \pi(L, R, O)$ for some **fixed permutation** π

Approach 1:

Permutation check: $f(x) = f(\pi(x)) \forall x \in H$

Approach 2:

 $\widetilde{\boldsymbol{\pi}} = mle(\boldsymbol{\pi})$



Advantage: Linear prover time Tradeoff: Not super fit for "small" fields

Advantage: soundness err = $(\log |H|)^2 / |\mathbb{F}|$ Tradeoff: quasilinear prover

Hyperplonk: Components



The last mile to success: batching



Outline

• Adapt Plonk to the Boolean hypercube

- High-degree gates support
- Hyperplonk+: support lookup on the Boolean hypercube
- Evaluations

Hyperplonk: High-degree Gates

$$S(y) \cdot [L(y) + R(y)] + (1 - S(y)) \cdot [L(y) \cdot R(y)] = O(y) \ \forall y \in H_{gates}$$

The gate formula can be more general
e.g. more selectors and more different custom gates

$$S(y) \cdot [L(y) + R(y)] + (1 - S(y)) \cdot [G(L(y), R(y))] = O(y) \ \forall y \in H_{gates}$$

Benefit: reduce the size of the circuit/witness -> faster prover

E.g.: ECC addition, Rescue hash, etc...

Hyperplonk: High-degree Gates

Plonk: more expensive quotient check

- Quotient polynomial q(X) has degree O(dn)
- Commit q(X) -> increase group operations
- Higher degree poly-mul -> larger-sized FFTs
- In practice, deg ≤ 8

Quotient Check $f(X) = Zero_H(X) \cdot \boldsymbol{q}(X)$

V.S.

Hyperplonk: efficient multivariate sumcheck• Only $\tilde{O}(dn)$ field-operations• Allows much higher degree	Sumcheck $\sum_{y \in B_{\mu}} f(y) \cdot eq \ (y,r) = 0$
---	--

Efficient Sumcheck for High-degree Polynomials

Goal: "convince" V that $\sum_{x \in B_{\mu}} f(X) = s$ VerifierProver $r_{\mu}(X) = \sum_{b \in B_{\mu-1}} f(b, X)$ Verifier $\alpha_u \leftarrow \mathbb{F}_p$ $check r_{\mu}(0) + r_{\mu}(1) = s$ \dots \dots $n_i \leftarrow \mathbb{F}_p$ $check r_i(0) + r_i(1) = r_{i+1}(\alpha_{i+1})$ \dots \dots </

Optimizations to the classic Sumcheck [LFKN92]

- Sending r_i as univariate oracles (d field elems -> 1 group elem)
- Decrease the number of queries per round: 3 queries -> 1 query
- More efficient algorithm for computing Sumcheck for high-degree polynomials
 - $O(2^u d^2) \rightarrow O(2^\mu d \log^2 d)$ (and $O(2^\mu d \log d)$ for certain custom gates)
 - In practice, replace FFTs with Karatsuba

No proof size/verifier dependence on d anymore

Evaluations: High-degree Gates



 Degree 32 is only 30% more expensive than degree 2

Outline

- Adapt Plonk to the Boolean hypercube
- High-degree gates support
- Hyperplonk+: support lookup on the Boolean hypercube
- Evaluations

Lookup Gates



The Idea of Plookup [GW20]



Shift in the Boolean Hypercube

Cool connection w/ short PCP Step 1: Shift on the Boolean hypercube from a multiplication in $GF(2^{\mu})$ [Ben-Sasson, Sudan 05] Let $p_{\mu}(X) \in \mathbb{F}_{2}^{\mu}[X]$ be a **primitive polynomial** $\mathbb{F}_{2}[X]/(p_{\mu}) \cong GF(2^{\mu})$ and X is a generator! E.g. $p_{\mu}(X) = X^4 + X + 1$ $b_2 X^4 + b_2 X^3 + b_1 X^2 + b_0 X$ $b_3 X^3 + b_2 X^2 + b_1 X + b_0$ multiply X $= b_2 X^3 + b_1 X^2 + (b_0 \oplus b_3) X + b_3$ $next(u) = (b_2, b_1, b_0 \oplus b_3, b_3) \in \mathbb{F}_{2^4} \setminus \{0\}$ $u = (b_3, b_2, b_1, b_0) \in \mathbb{F}_{2^4} \setminus \{0\}$ 0001 0100 Caveats 0010 1000 0011 0110 $b_0 \oplus b_3 = b_0 + b_3 - 2b_0b_3$ is quadratic 1110 0111-1010 0101 1011 1100 Individual deg of f(next(X)) = 2• Naïve MSET check on f(next(X))• 1101 1001 1111 • $\Theta(n^2)$ sumcheck prover

Shift in the Boolean Hypercube

Step 2: Simulate f(next(X)) on the Boolean hypercube using the multilinear f(X)!

 $u = (b_3, b_2, b_1, b_0) \in \mathbb{F}_{2^4} \setminus \{0\} \xrightarrow{\text{multiply X}} \text{next}(u) = (b_2, b_1, b_0 \oplus b_3, b_3) \in \mathbb{F}_{2^4} \setminus \{0\}$

Define

Idea: enumerate b_3 and linearize next()!

$$f_{\mathcal{O}}(X_3, \dots, X_0) \coloneqq X_3 \cdot f(X_2, X_1, 1 - X_0, 1) + (1 - X_3) \cdot f(X_2, X_1, X_0, 0)$$

Property 1: $f_{U}(X) = f(next(X))$ on the Boolean hypercube!

$$\begin{aligned} \text{Case } b_3 &= 1: \\ f_{\bigcup}(b_3, \dots, b_0) &\coloneqq 1 \cdot f(b_2, b_1, 1 - b_0, 1) &+ (1 - 1) \\ &= f(b_2, b_1, b_0 \oplus b_3, b_3) \\ &= f(\text{next}(b_3, \dots, b_0)) \\ \text{Case } b_3 &= 0: \\ f_{\bigcup}(b_3, \dots, b_0) &\coloneqq 0 \cdot &+ (1 - 0) \cdot f(b_2, b_1, b_0, 0) = f(b_2, b_1, b_0 \oplus b_3, b_3) \\ &= f(\text{next}(b_3, \dots, b_0)) \end{aligned}$$

Shift in the Boolean Hypercube

Step 2: Simulate f(next(X)) on the Boolean hypercube using the multilinear f(X)! $u = (b_3, b_2, b_1, b_0) \in \mathbb{F}_{2^4} \setminus \{0\} \xrightarrow{\text{multiply X}} \text{next}(u) = (b_2, b_1, b_0 \oplus b_3, b_3) \in \mathbb{F}_{2^4} \setminus \{0\}$ **Idea**: enumerate b_3 and linearize next()! Define $f_{\bigcup}(X_3, \dots, X_0) \coloneqq X_3 \quad f(X_2, X_1, 1 - X_0, 1) + (1 - X_3) \quad f(X_2, X_1, X_0, 0)$ **Property 2:** Can simulate eval of f_{\bigcirc} from 2 evals of fNo need to build/check The lookup PolyIOP can run sumcheck given oracles T, hT(next(X)), h(next(X))lookup PIOP with linear prover time!

Outline

- Adapt Plonk to the Boolean hypercube
- High-degree gates support
- Hyperplonk+: support lookup on the Boolean hypercube
- Evaluations

Evaluations

	2^{10}	2^{11}	2^{12}	2^{13}	2^{14}	2^{15}	2^{16}	2^{17}	2^{18}	2^{19}	2^{20}
HP (Vanilla Gate)	0.07	0.1	0.14	0.2	0.3	0.5	0.8	1.4	2.5	5.1	9.6
HP (Jellyfish Gate)	0.1	0.13	0.18	0.27	0.4	0.67	1.2	2	3.7	7.3	13.5
Jellyfish Plonk	0.07	0.1	0.15	0.25	0.46	0.78	1.4	2.7	5.5	10.8	22
Ark-Spartan	0.51	0.72	0.9	1.4	1.9	3.1	4.7	8.3	13.7	27	44

Table 6: 64-thread prover's performance (in seconds) for varying number of constraints under different schemes.

Application	$\mathcal{R}_{\mathrm{R1}CS}$	Ark-Spartan	$\mathcal{R}_{\mathrm{PLONK+}}$	Jellyfish	HyperPlonk
3-to-1 Rescue Hash	288 [1]	$422 \mathrm{\ ms}$	144 [71]	$40 \mathrm{ms}$	$88 \mathrm{ms}$
PoK of Exponent	3315 [63]	$902 \mathrm{\ ms}$	783 [63]	$64 \mathrm{ms}$	$105 \mathrm{\ ms}$
ZCash circuit	2^{17} [55]	8.3 s	2^{15} [42]	0.8 s	0.6 s
Zexe's recursive circuit	2^{22} [81]	$6 \min$	2^{17} [81]	13.1 s	$5.1 \mathrm{~s}$
Rollup of 50 private tx	2^{25}	$39 \min^b$	2^{20} [71]	110 s	38.2 s
$z k EVM circuit^{a}$	N/A	N/A	2^{27}	$1 \text{ hour}^{b,c}$	$25 \min^{b,c}$

Proof size (μ =20)(with KZG PCS)

- Plonk ~1kb
- HP 4.7kb
- Spartan 40kb

Implementation

- Jellyfish: highly optimized
- HP: many possible improvements
- Prover time
 - HP outperforms Plonk at 2¹⁴; at 2²⁰ HP is 60% faster then Plonk
 - Fewer constraints needed than Spartan for the same application



Cost Breakdown

Bottleneck

PCS commitments and batch openings



Potential optimizations

- many eval points and polys are identical
- Batching all zero-checks

Further improvements on hardware

• Multi-exps and sumchecks are highly parallelizable and hardware-friendly

Summary & Open problems

HyperPlonk+ PolyIOP

- Linear-time prover, no use of FFT-friendly field
- Native high-degree custom gate support
- Lookup gate support with linear-time prover
- Concretely efficient

Open problems

- Degree k > 1 shifts?
 - Our technique: needs 2^k evals on f to simulate an eval on $f_{\mathcal{O}}$
- Linear-time permutation argument for fields with small characteristics

Thank you!

https://eprint.iacr.org/2022/1355.pdf

Permutation Check

Prove $(L, R, O) = \pi(L, R, O)$ for some **fixed permutation** π



Permutation Check

Prove $(L, R, O) = \pi(L, R, O)$ for some **fixed permutation** π



Efficient Batch Evaluations

Goal: "convince" V that $f_i(z_i) = y_i \forall i \in [k] (f_i \in \mathbb{F}_p[X_1, ..., X_\mu])$

Prover $P(S_p, \boldsymbol{x}, \boldsymbol{w})$

PolyIOP for batch evaluate

Verifier V(S_v, x)

A **single** evaluation query on a multilinear polynomial

Requirements

- Prover time \approx the time complexity of k poly evaluations
- Logarithmic communication/rounds and verification time

Our results

- Prover time $O(k2^{\mu})$ (optimal)
- $k\mu$ -factor improvement (k>13 in HP)
- $O(\mu)$ communication
- $O(k\mu)$ verification time

State-of-the-arts

- Halo-infinite [BDFG20] only for univariate polys
- Thaler's solution has prover time $O(k^2 \mu \cdot 2^{\mu})$



The Idea of Plookup [GW20]



HyperPlonk+: Open Problems

- Shifting by k > 1 is useful in some applications (e.g. Halo2)
- Shifting by k is not as easy as shifting by 1
 - The linearization trick has an exponential blow-up on the shift distance
 - Distance-k shift: 2^k evals on f to simulate an eval on f_{\cup}
- Can we implement a distance-k shift more efficiently?

