

# Lecture 6

# Recap

- Two lectures ago: Our first application of the sum-check protocol.
  - An IP for #SAT with a polynomial-time verifier.
  - **P** ran in time exponential in the input size.
  - But the fastest known algorithm for this problem requires exponential time.
    - So can't really hope for a faster prover for this problem.
- Last lecture we saw some doubly-efficient IPs.
  - **V** runs in **linear** time.
  - **P** runs in polynomial time.
    - In fact, we achieved “super-efficiency”.
    - meaning **P** ran the fastest known algorithm for the problem, and then did a low-order amount of additional work to prove correctness.
    - Counting triangles, matrix multiplication.

Today: A **General-Purpose**  
Doubly-Efficient Interactive  
Proof

# General-Purpose Interactive Proof and Argument Implementations

- Start with a computer program written in high-level programming language (C, Java, etc.)
- Step 1: Turn the program into an equivalent model amenable to probabilistic checking.
  - Typically some type of arithmetic circuit.
  - Called the **Front End** of the system.
- Step 2: Run an interactive proof or argument on the circuit.
  - Called the **Back End** of the system.

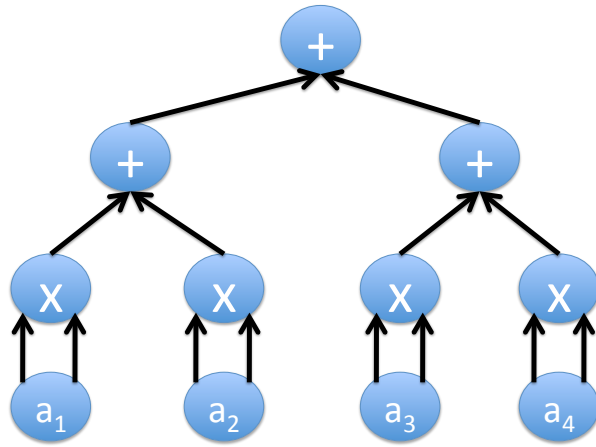
```
blink.c
.....
* Author: Leehi Szechtel
* Filename: blink.c
* Chip: at32mc12
*/

#define F_CPU 1000000
#include <avr/io.h>
#include <avr/delay.h>
#include "../leah_library/pin_waves.h"

int main(void)
{
  b0output();
  b1input();
  b1high();

  for(;;)
  {
    if (b1isLow())
    {
      b0high();
    }
    else
    {
      b0Low();
    }
  }
  return 0;
}
```

Front End



P and V run interactive proof or argument system (back end) on circuit

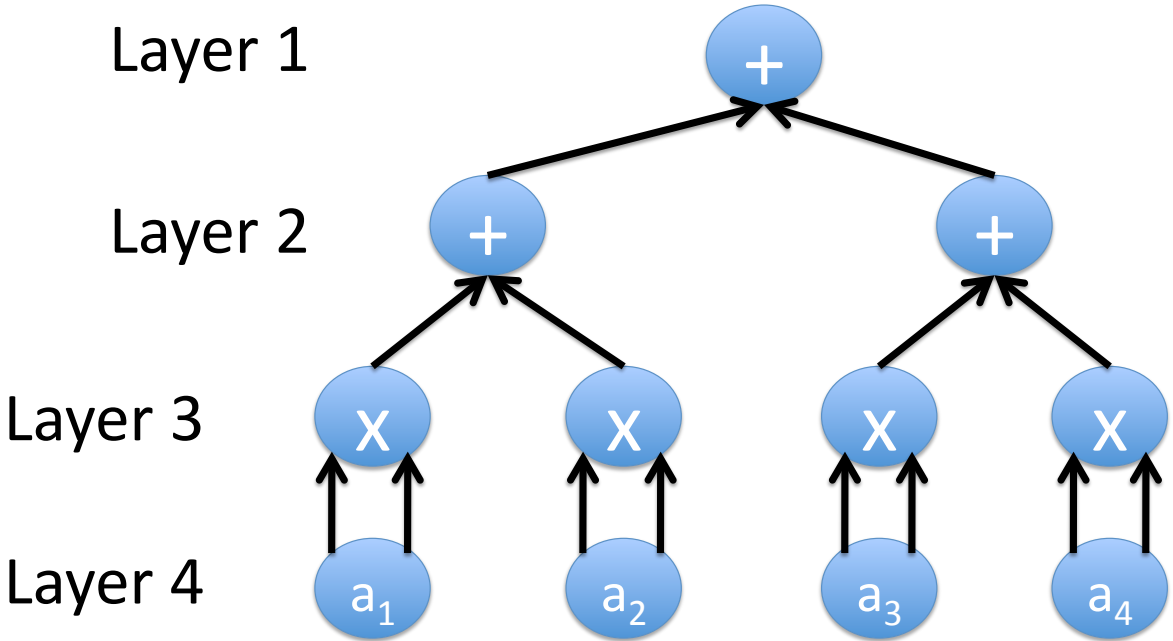
# Sources of Prover Overhead in VC Systems

Source of Overhead	<b>P</b> Overhead vs. Native <b>(Crude Estimate)</b>	Slowdown Depends On...
Front End (overhead due to using a circuit representation of the computation)	(ratio of circuit size to number of machine steps of original program) <b>1x-10,000x</b>	<ul style="list-style-type: none"><li>• How amenable is the high-level computer program is to representation via circuits?</li><li>• What type of circuits can the back-end handle?</li></ul>
Back-End	(ratio of <b>P</b> time to evaluating circuit gate-by-gate) <b>10x-1,000x</b>	<ul style="list-style-type: none"><li>• Varies by back-end and computation structure (e.g., data parallel?)</li></ul>

# The GKR Protocol

A General-Purpose Doubly-Efficient  
Interactive Proof

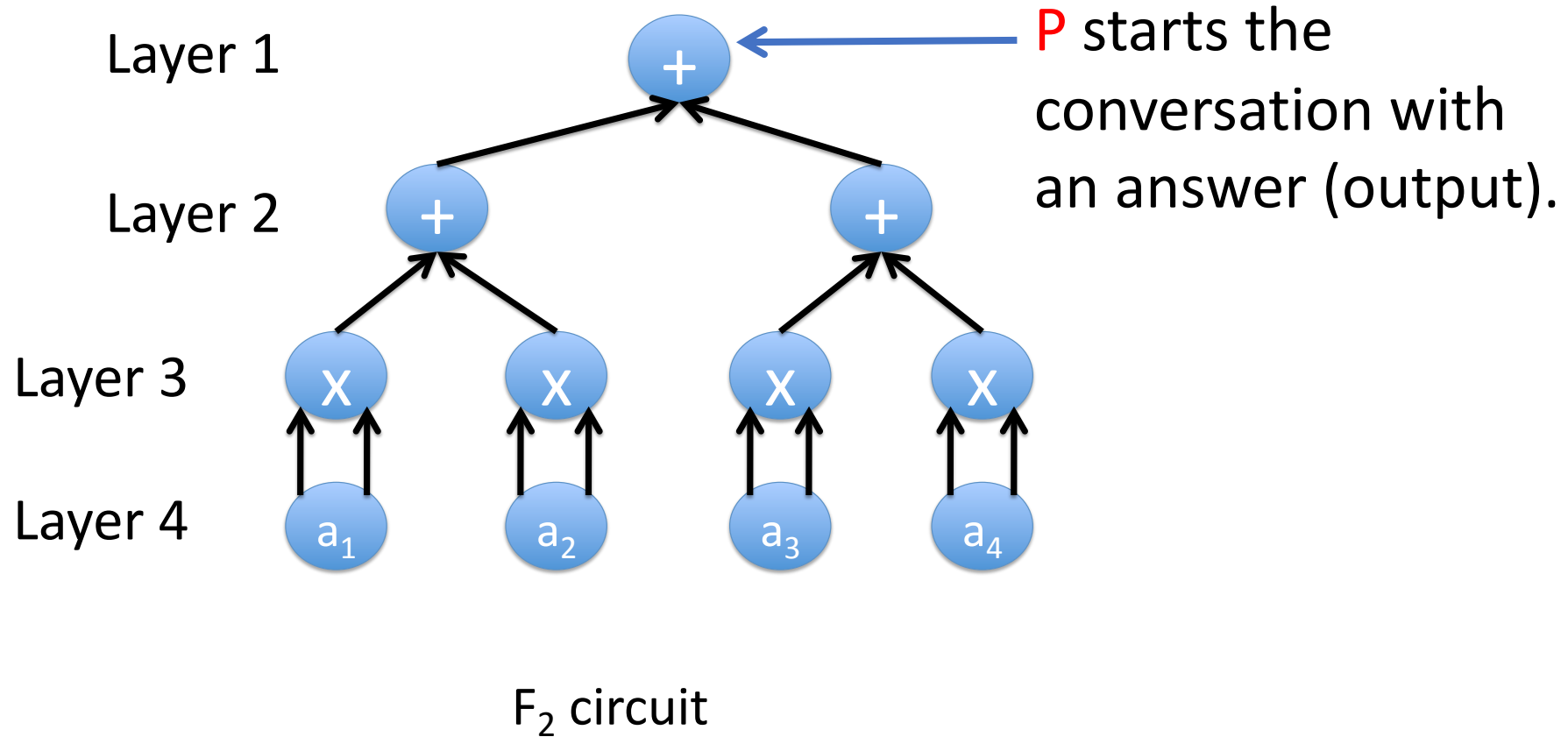
# The GKR Protocol: Overview



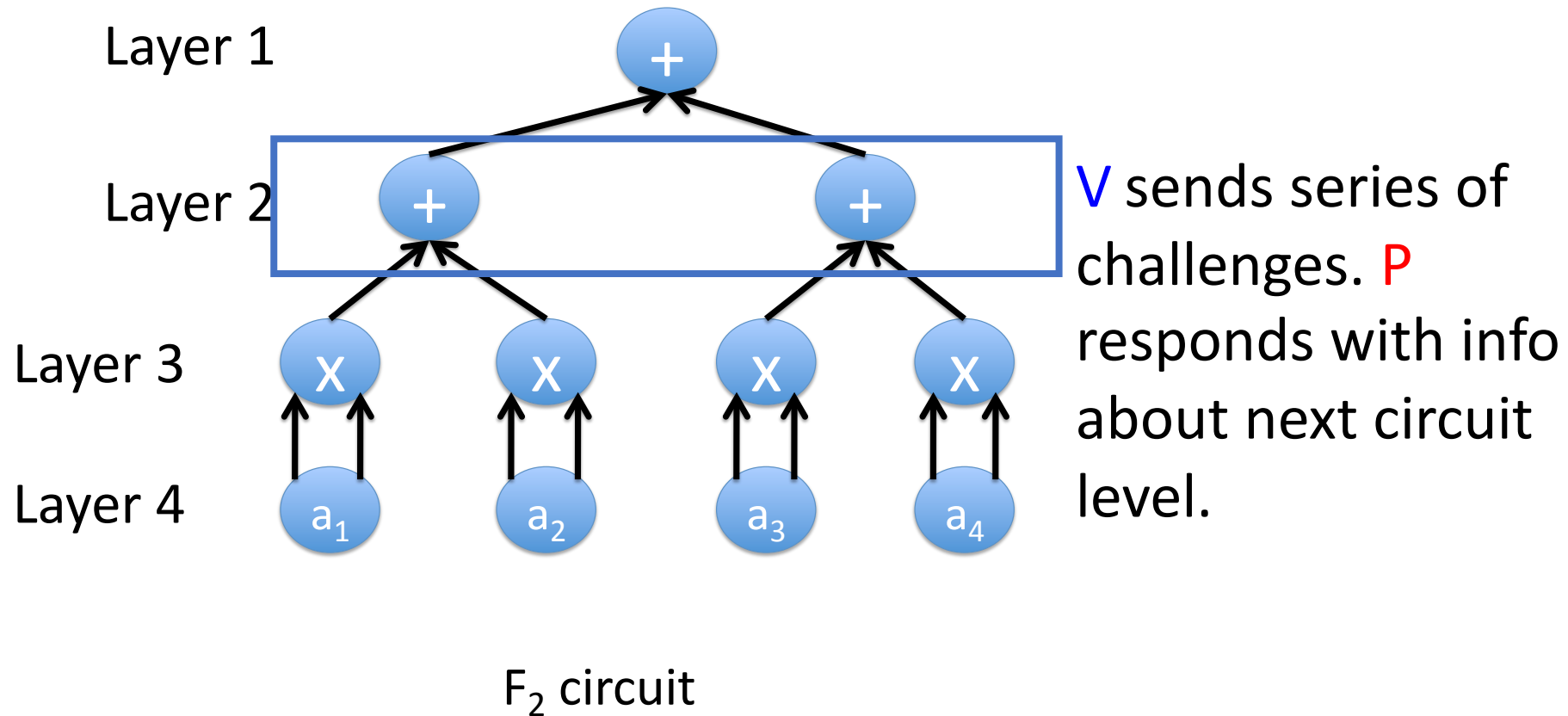
$F_2$  circuit



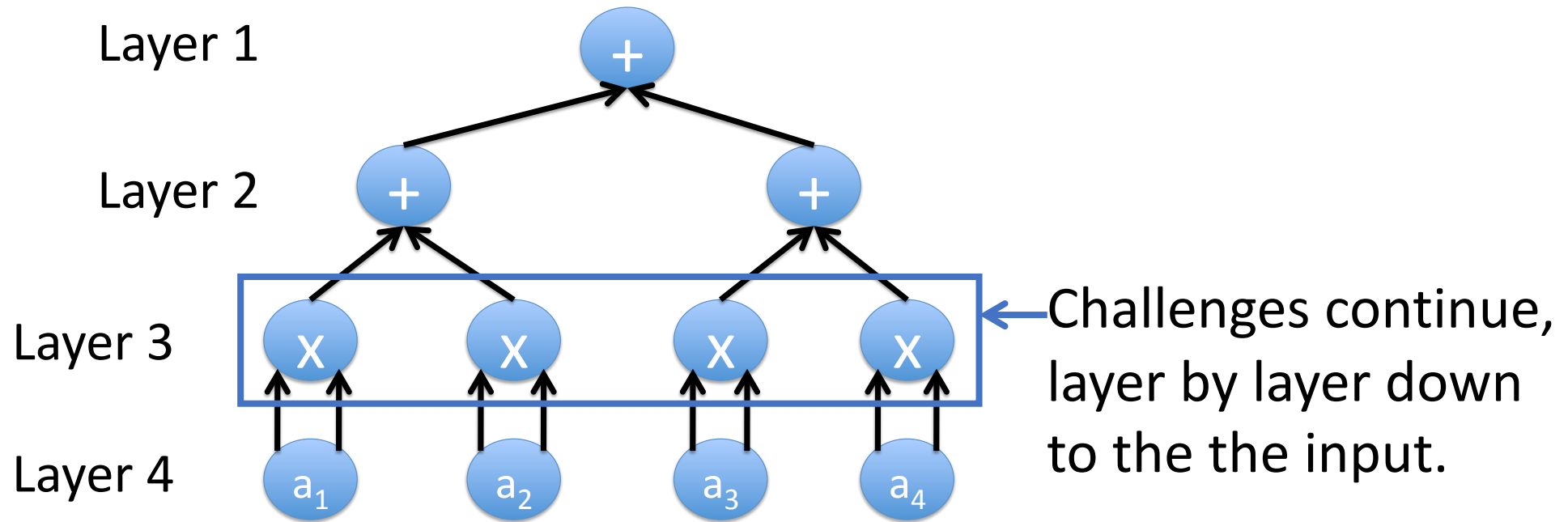
# The GKR Protocol: Overview



# The GKR Protocol: Overview

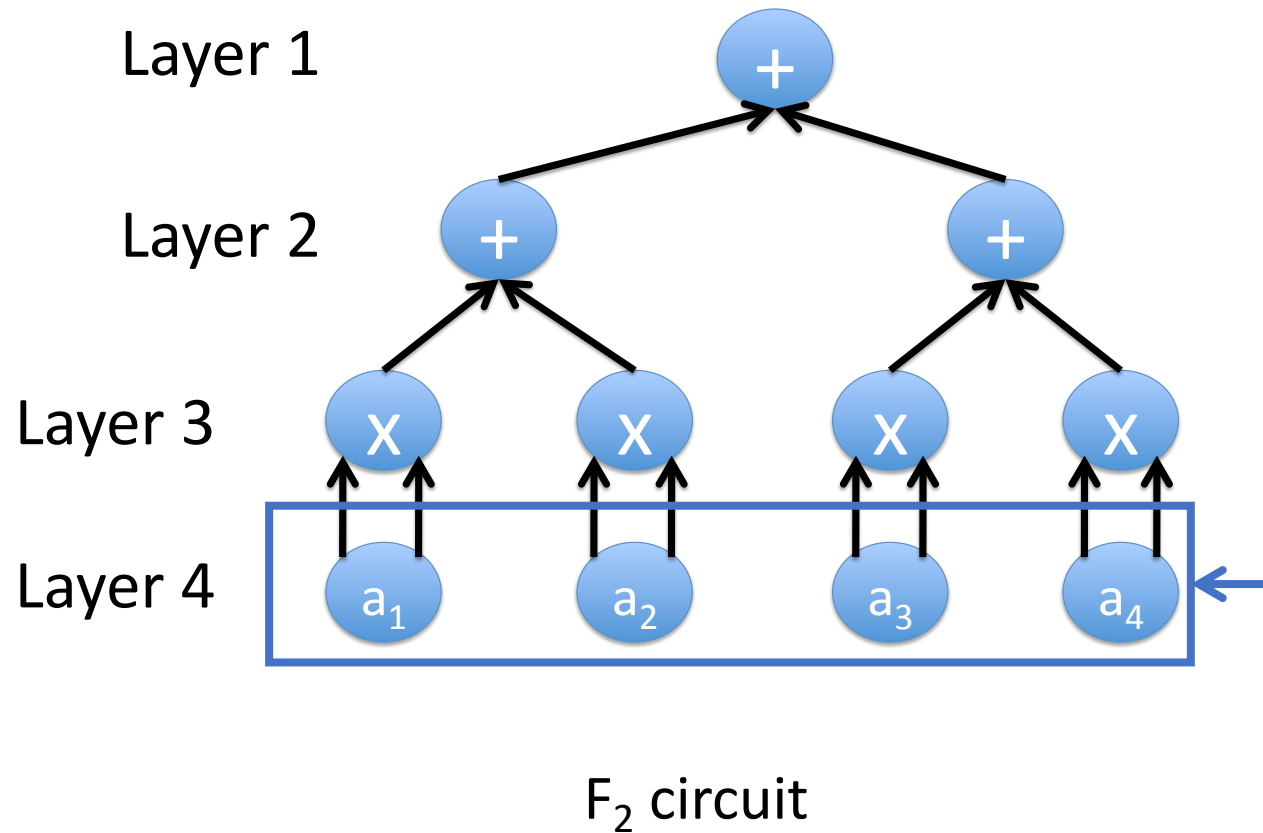


# The GKR Protocol: Overview



$F_2$  circuit

# The GKR Protocol: Overview



Finally, **P** says something about the (multilinear extension of the) input.

# Notation

- Assume layers  $i$  and  $i + 1$  of  $C$  have  $S$  gates each.
  - Assign each gate a binary label ( $\log S$  bits).
- Let  $W_i(\mathbf{a}): \{0,1\}^{\log S} \rightarrow \mathbf{F}$  output the value of gate  $\mathbf{a}$  at layer  $i$ .

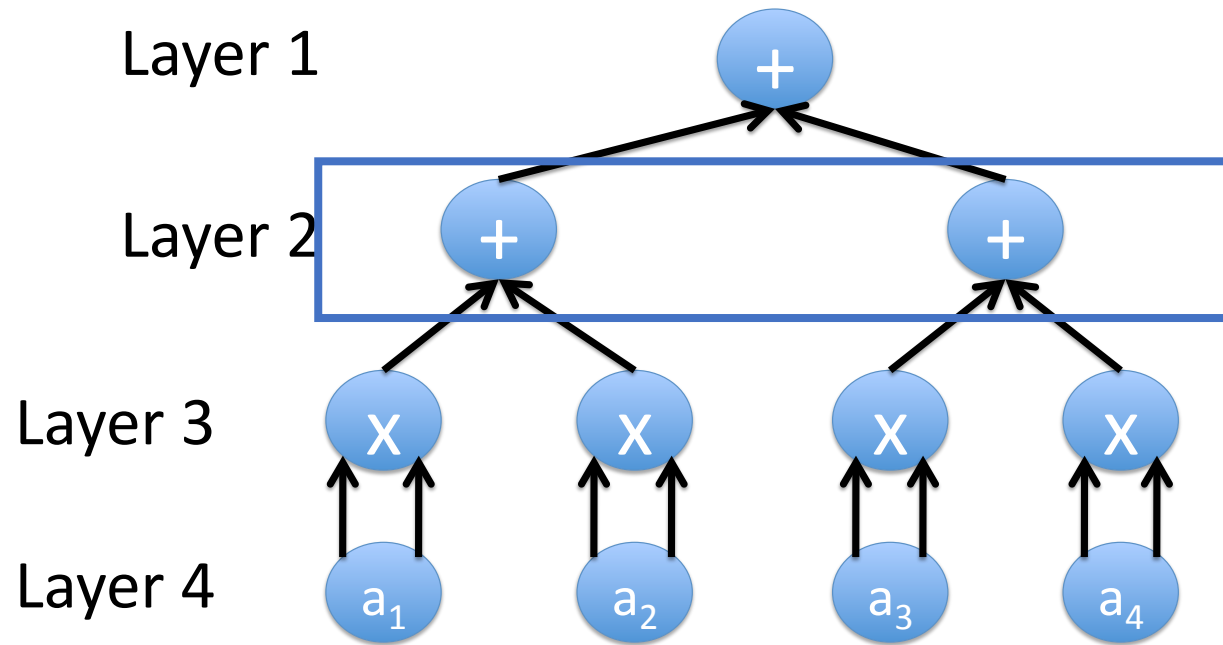
# Notation

- Assume layers  $i$  and  $i + 1$  of  $C$  have  $S$  gates each.
  - Assign each gate a binary label ( $\log S$  bits).
- Let  $W_i(\mathbf{a}): \{0,1\}^{\log S} \rightarrow \mathbf{F}$  output the value of gate  $\mathbf{a}$  at layer  $i$ .
- Let  $\text{add}_i(\mathbf{a}, \mathbf{b}, \mathbf{c}): \{0,1\}^{3 \log S} \rightarrow \mathbf{F}$  output 1 iff  $(\mathbf{b}, \mathbf{c}) = (\text{in}_1(\mathbf{a}), \text{in}_2(\mathbf{a}))$  and gate  $\mathbf{a}$  at layer  $i$  is an addition gate.

# Notation

- Assume layers  $i$  and  $i + 1$  of  $C$  have  $S$  gates each.
  - Assign each gate a binary label ( $\log S$  bits).
- Let  $W_i(\mathbf{a}): \{0,1\}^{\log S} \rightarrow \mathbf{F}$  output the value of gate  $\mathbf{a}$  at layer  $i$ .
- Let  $\text{add}_i(\mathbf{a}, \mathbf{b}, \mathbf{c}): \{0,1\}^{3 \log S} \rightarrow \mathbf{F}$  output 1 iff  $(\mathbf{b}, \mathbf{c}) = (\text{in}_1(\mathbf{a}), \text{in}_2(\mathbf{a}))$  and gate  $\mathbf{a}$  at layer  $i$  is an addition gate.
- Let  $\text{mult}_i(\mathbf{a}, \mathbf{b}, \mathbf{c}): \{0,1\}^{3 \log S} \rightarrow \mathbf{F}$  output 1 iff  $(\mathbf{b}, \mathbf{c}) = (\text{in}_1(\mathbf{a}), \text{in}_2(\mathbf{a}))$  and gate  $\mathbf{a}$  at layer  $i$  is a multiplication gate.

# The GKR Protocol: Overview



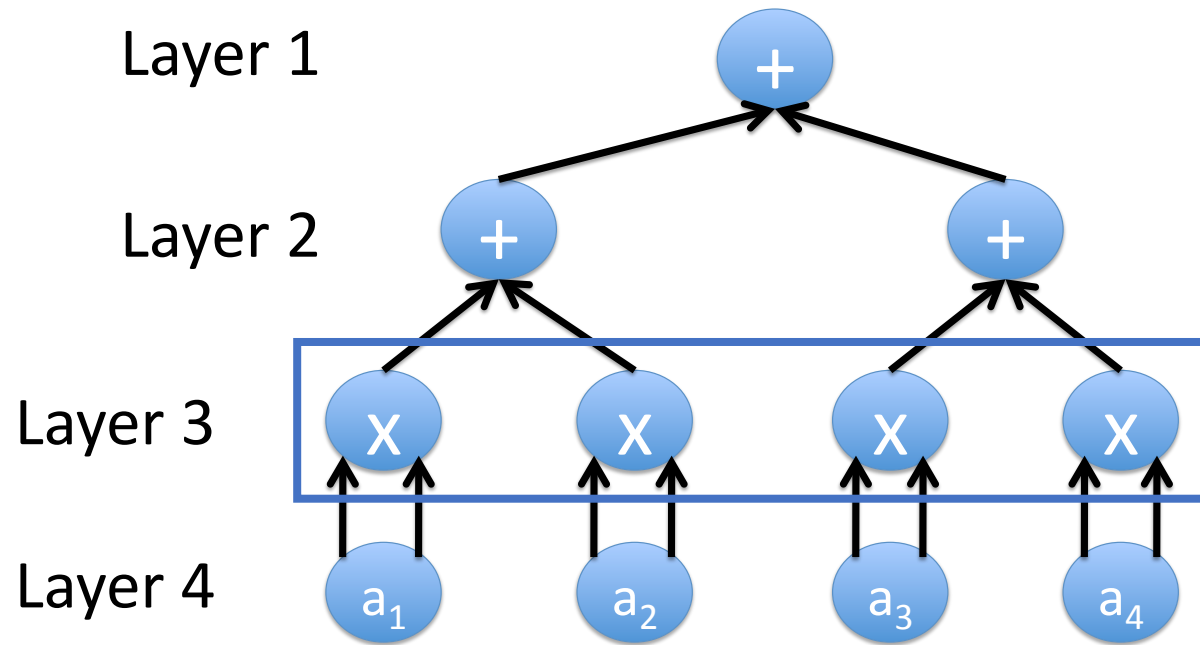
$$\text{add}_2(0, (0, 0), (0, 1)) = 1$$

$$\text{add}_2(1, (1, 0), (1, 1)) = 1$$

$F_2$  circuit



# The GKR Protocol: Overview



$F_2$  circuit

$$\text{mult}_3((0,0), (0,0), (0,0)) = 1$$

$$\text{mult}_3((0,1), (0,1), (0,1)) = 1$$

$$\text{mult}_3((1,0), (1,0), (1,0)) = 1$$

$$\text{mult}_3((1,1), (1,1), (1,1)) = 1$$

# GKR Protocol: Goal of Iteration $i$

- Iteration  $i$  starts with a claim from  $\mathbf{P}$  about  $\tilde{W}_i(\mathbf{r}_1)$  for a random point  $\mathbf{r}_1 \in \mathbf{F}^{\log S}$ .
- Goal: Reduce this to a claim about  $\tilde{W}_{i+1}(\mathbf{r}_2)$  for a random point  $\mathbf{r}_2 \in \mathbf{F}^{\log S}$ .
- Observation:  $W_i(\mathbf{a}) =$

$$\sum_{\mathbf{b}, \mathbf{c} \in \{0,1\}^{\log S}} [\text{add}_i(\mathbf{a}, \mathbf{b}, \mathbf{c})(W_{i+1}(\mathbf{b}) + W_{i+1}(\mathbf{c})) + \text{mult}_i(\mathbf{a}, \mathbf{b}, \mathbf{c})(W_{i+1}(\mathbf{b}) \cdot W_{i+1}(\mathbf{c}))]$$

- Hence, the following equality holds as formal polynomials:

$$\tilde{W}_i(\mathbf{a}) = \sum_{\mathbf{b}, \mathbf{c} \in \{0,1\}^{\log S}} [\widetilde{\text{add}}_i(\mathbf{a}, \mathbf{b}, \mathbf{c})(\tilde{W}_{i+1}(\mathbf{b}) + \tilde{W}_{i+1}(\mathbf{c})) + \widetilde{\text{mult}}_i(\mathbf{a}, \mathbf{b}, \mathbf{c})(\tilde{W}_{i+1}(\mathbf{b}) \cdot \tilde{W}_{i+1}(\mathbf{c}))]$$

# GKR Protocol: Goal of Iteration $i$

- So  $V$  applies sum-check protocol to compute

- $\tilde{W}_i(\mathbf{r}_1) = \sum_{\mathbf{b}, \mathbf{c} \in \{0,1\}^{\log s}} g(\mathbf{b}, \mathbf{c})$ , where:

$$g(\mathbf{b}, \mathbf{c}) = \widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\tilde{W}_{i+1}(\mathbf{b}) + \tilde{W}_{i+1}(\mathbf{c})) \\ + \widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\tilde{W}_{i+1}(\mathbf{b}) \cdot \tilde{W}_{i+1}(\mathbf{c}))$$

# GKR Protocol: Goal of Iteration $i$

- So  $V$  applies sum-check protocol to compute

- $\tilde{W}_i(\mathbf{r}_1) = \sum_{\mathbf{b}, \mathbf{c} \in \{0,1\}^{\log s}} g(\mathbf{b}, \mathbf{c})$ , where:

$$g(\mathbf{b}, \mathbf{c}) = \widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\tilde{W}_{i+1}(\mathbf{b}) + \tilde{W}_{i+1}(\mathbf{c})) \\ + \widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\tilde{W}_{i+1}(\mathbf{b}) \cdot \tilde{W}_{i+1}(\mathbf{c}))$$

- At end of sum-check protocol,  $V$  must evaluate  $g(\mathbf{r}_2, \mathbf{r}_3)$ .

# GKR Protocol: Goal of Iteration $i$

- So  $V$  applies sum-check protocol to compute

- $\tilde{W}_i(\mathbf{r}_1) = \sum_{\mathbf{b}, \mathbf{c} \in \{0,1\}^{\log s}} g(\mathbf{b}, \mathbf{c})$ , where:

$$g(\mathbf{b}, \mathbf{c}) = \widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\tilde{W}_{i+1}(\mathbf{b}) + \tilde{W}_{i+1}(\mathbf{c})) \\ + \widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\tilde{W}_{i+1}(\mathbf{b}) \cdot \tilde{W}_{i+1}(\mathbf{c}))$$

- At end of sum-check protocol,  $V$  must evaluate  $g(\mathbf{r}_2, \mathbf{r}_3)$ .
- Let us assume  $V$  can compute  $\widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3)$  and  $\widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3)$  unaided in time  $\text{polylog}(n)$ .
- Then  $V$  only needs to know  $\tilde{W}_{i+1}(\mathbf{r}_2)$  and  $\tilde{W}_{i+1}(\mathbf{r}_3)$  to complete this check.

# GKR Protocol: Goal of Iteration $i$

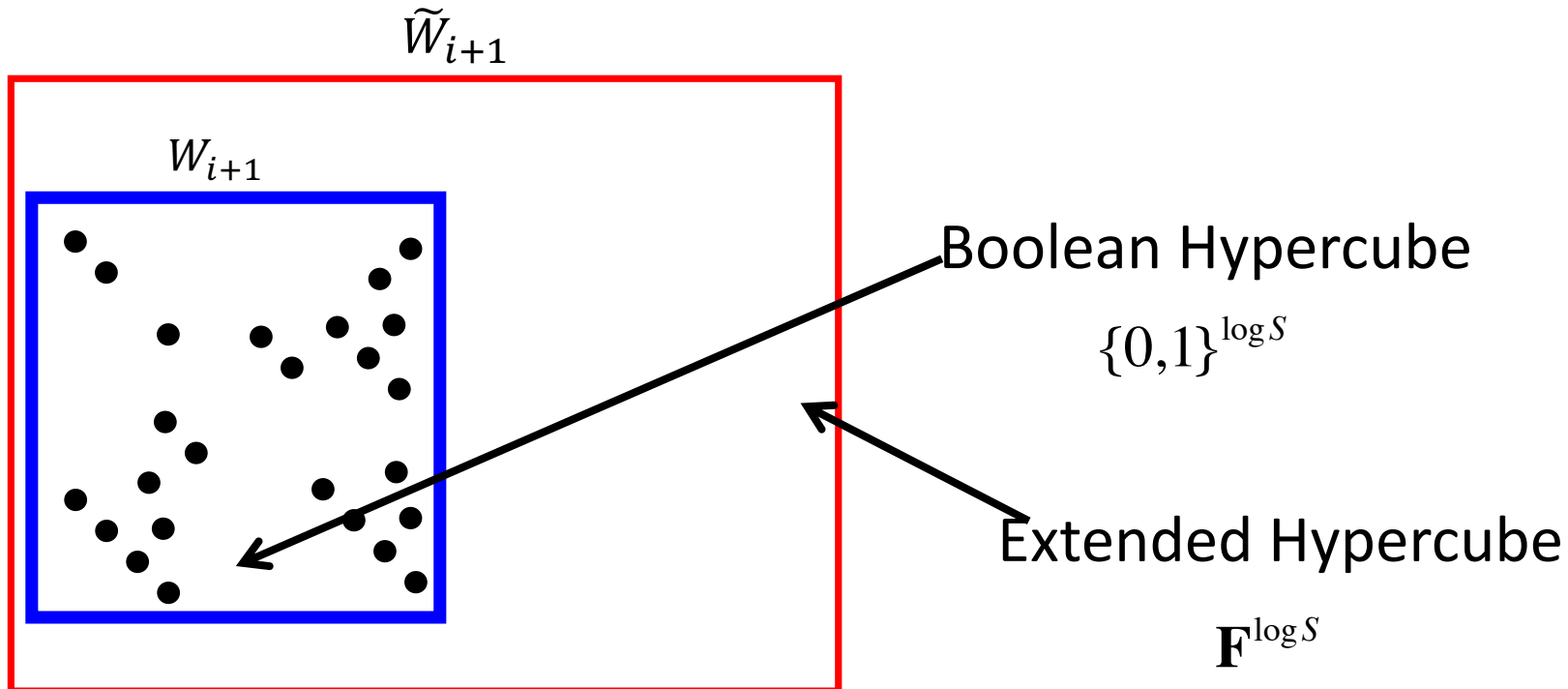
- So  $V$  applies sum-check protocol to compute
- $\tilde{W}_i(\mathbf{r}_1) = \sum_{\mathbf{b}, \mathbf{c} \in \{0,1\}^{\log s}} g(\mathbf{b}, \mathbf{c})$ , where:
$$g(\mathbf{b}, \mathbf{c}) = \widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\tilde{W}_{i+1}(\mathbf{b}) + \tilde{W}_{i+1}(\mathbf{c})) \\ + \widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\tilde{W}_{i+1}(\mathbf{b}) \cdot \tilde{W}_{i+1}(\mathbf{c}))$$
- At end of sum-check protocol,  $V$  must evaluate  $g(\mathbf{r}_2, \mathbf{r}_3)$ .
- Let us assume  $V$  can compute  $\widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3)$  and  $\widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3)$  unaided in time  $\text{polylog}(n)$ .
- Then  $V$  only needs to know  $\tilde{W}_{i+1}(\mathbf{r}_2)$  and  $\tilde{W}_{i+1}(\mathbf{r}_3)$  to complete this check.
- Iteration  $i + 1$  is devoted to computing these values.

# Remaining Issue: Reducing to Verification of a Single Point

- There is one remaining problem: we don't want to have to separately verify both  $\tilde{W}_{i+1}(\mathbf{r}_2)$  and  $\tilde{W}_{i+1}(\mathbf{r}_3)$  in iteration  $i + 1$ .
- Solution: Reduce verifying both of the above values to verifying  $\tilde{W}_{i+1}(\mathbf{r}_4)$  for a single point  $\mathbf{r}_4$ .

# Remaining Issue: Reducing to Verification of a Single Point

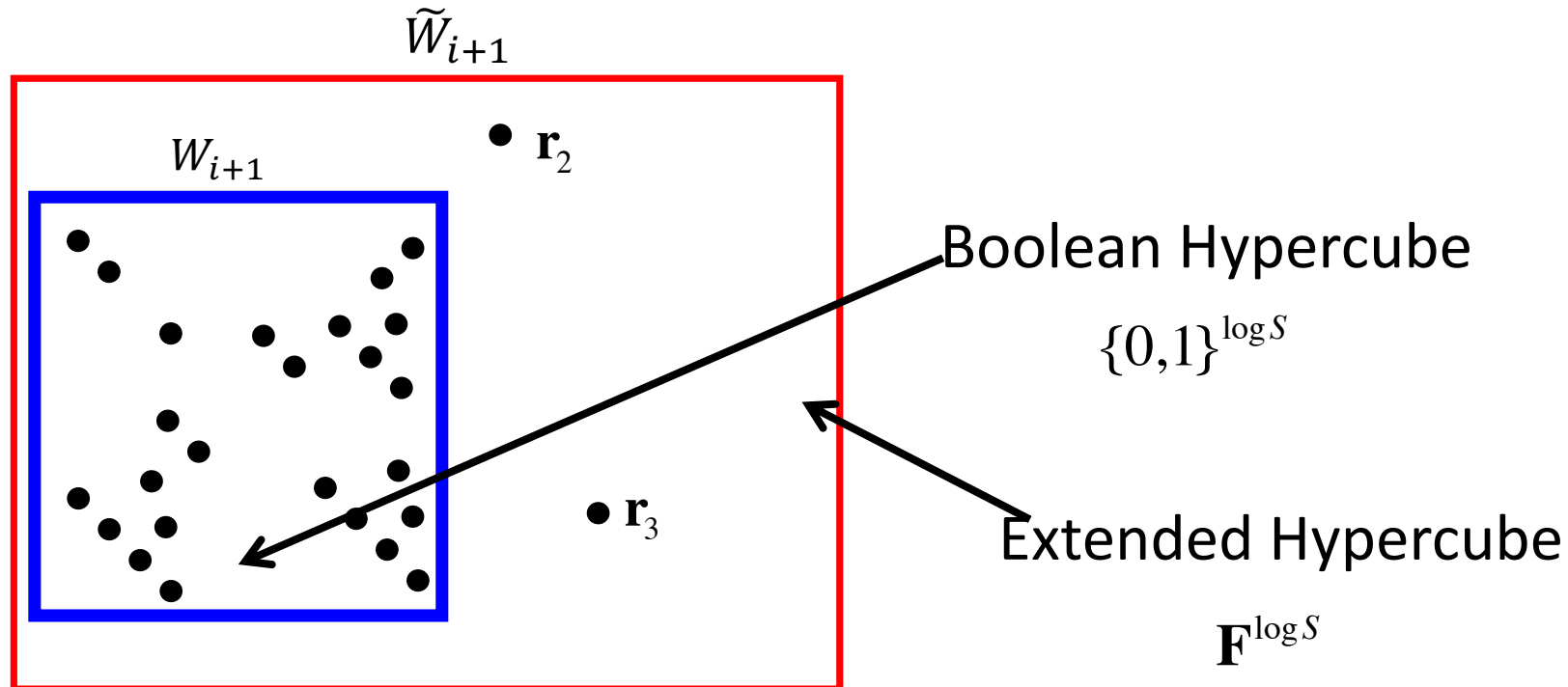
- There is one remaining problem: we don't want to have to separately verify both  $\tilde{W}_{i+1}(\mathbf{r}_2)$  and  $\tilde{W}_{i+1}(\mathbf{r}_3)$  in iteration  $i + 1$ .
- Solution: Reduce verifying both of the above values to verifying  $\tilde{W}_{i+1}(\mathbf{r}_4)$  for a single point  $\mathbf{r}_4 \in \mathbf{F}^{\log S}$ .





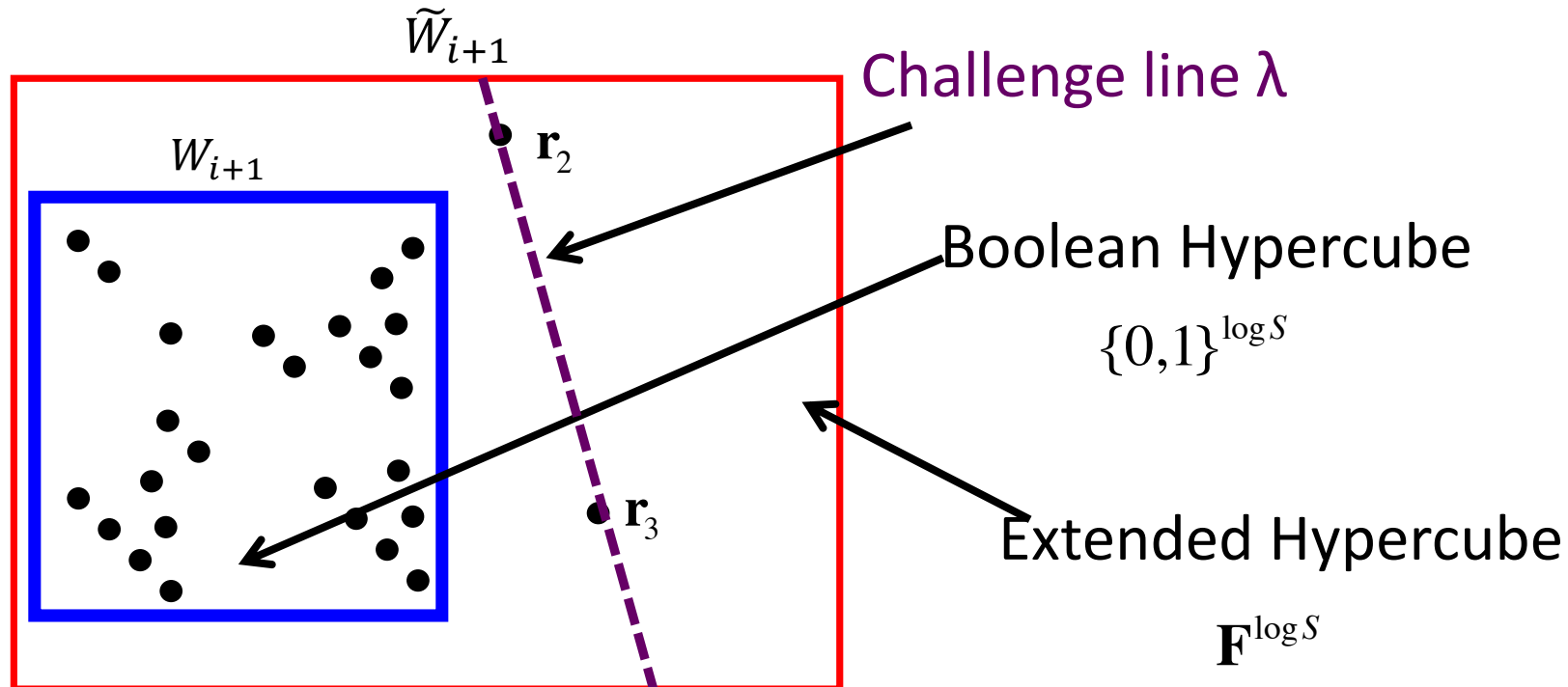
# Remaining Issue: Reducing to Verification of a Single Point

- There is one remaining problem: we don't want to have to separately verify both  $\tilde{W}_{i+1}(\mathbf{r}_2)$  and  $\tilde{W}_{i+1}(\mathbf{r}_3)$  in iteration  $i+1$ .
- Solution: Reduce verifying both of the above values to verifying  $\tilde{W}_{i+1}(\mathbf{r}_4)$  for a single point  $\mathbf{r}_4 \in \mathbf{F}^{\log S}$ .



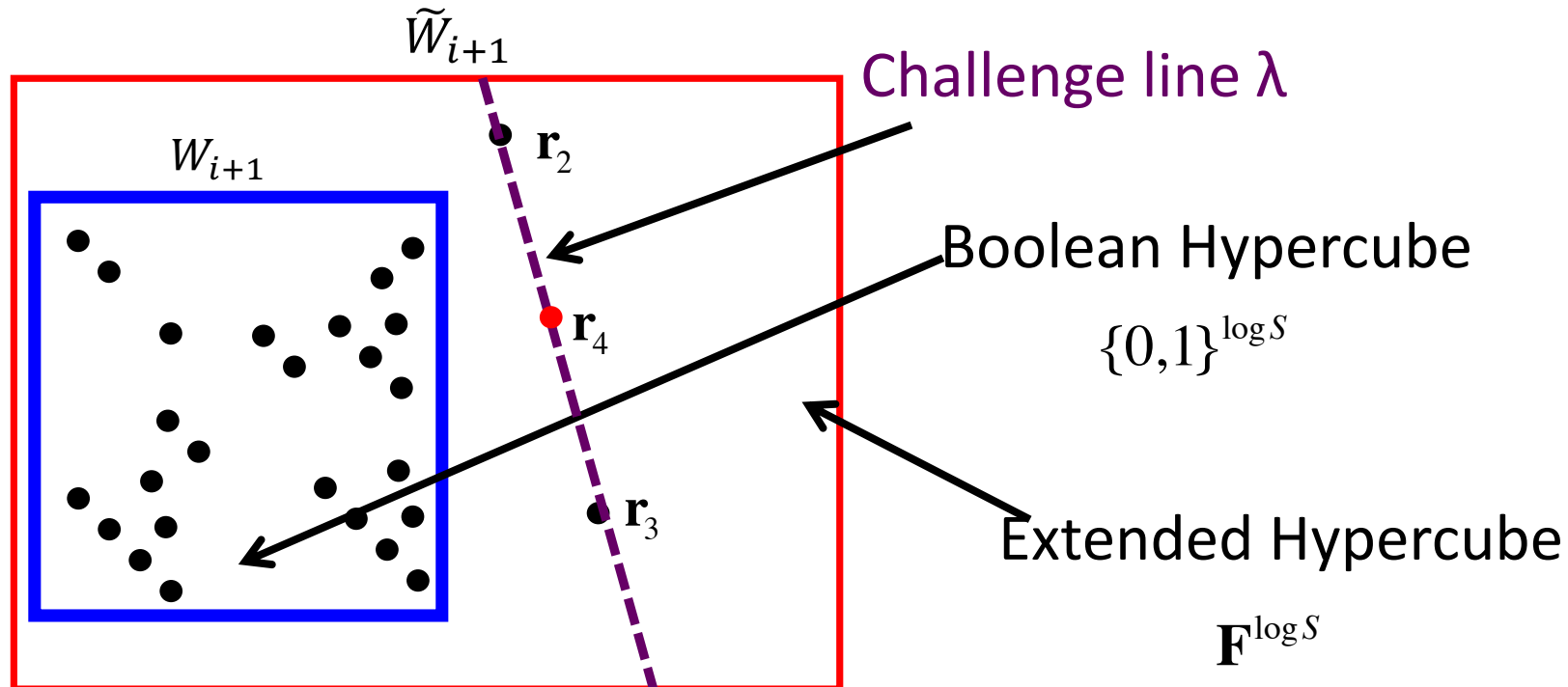
# Remaining Issue: Reducing to Verification of a Single Point

- There is one remaining problem: we don't want to have to separately verify both  $\tilde{W}_{i+1}(\mathbf{r}_2)$  and  $\tilde{W}_{i+1}(\mathbf{r}_3)$  in iteration  $i+1$ .
- Solution: Reduce verifying both of the above values to verifying  $\tilde{W}_{i+1}(\mathbf{r}_4)$  for a single point  $\mathbf{r}_4 \in \mathbf{F}^{\log S}$ .



# Remaining Issue: Reducing to Verification of a Single Point

- There is one remaining problem: we don't want to have to separately verify both  $\tilde{W}_{i+1}(\mathbf{r}_2)$  and  $\tilde{W}_{i+1}(\mathbf{r}_3)$  in iteration  $i+1$ .
- Solution: Reduce verifying both of the above values to verifying  $\tilde{W}_{i+1}(\mathbf{r}_4)$  for a single point  $\mathbf{r}_4 \in \mathbf{F}^{\log S}$ .



# Costs of the GKR protocol

- $V$  time is  $O(n + D \log S)$  where  $n$  is input size,  $D$  is circuit depth, and  $S$  is circuit size.
  - Assumes  $V$  can compute  $\widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3)$  and  $\widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3)$  unaided in time  $\text{polylog}(n)$
- Communication cost is  $O(D \log S)$ .



# Costs of the GKR protocol

- **V** time is  $O(n + D \log S)$  where  $n$  is input size,  $D$  is circuit depth, and  $S$  is circuit size.
  - Assumes **V** can compute  $\widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3)$  and  $\widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3)$  unaided in time  $\text{polylog}(n)$
- Communication cost is  $O(D \log S)$ .
- **P** time is  $O(S)$ .
  - A naïve implementation of **P** takes  $\Omega(S^3)$  time, where  $S$  is circuit size.
  - A sequence of works has brought this down to  $O(S)$ , for arbitrary circuits [CMT12, Thaler13, WJBSTWW17, XZZPS19]



# GKR Prover Runtime: Details

Recall: Core of the GKR protocol is applying sum-check to compute  $\sum_{\mathbf{b}, \mathbf{c} \in \{0,1\}^{\log s}} g(\mathbf{b}, \mathbf{c})$  where

$$g(\mathbf{b}, \mathbf{c}) = \widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\widetilde{W}_{i+1}(\mathbf{b}) + \widetilde{W}_{i+1}(\mathbf{c})) \\ + \widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\widetilde{W}_{i+1}(\mathbf{b}) \cdot \widetilde{W}_{i+1}(\mathbf{c}))$$

# GKR Prover Runtime: Details

Recall: Core of the GKR protocol is applying sum-check to compute  $\sum_{\mathbf{b}, \mathbf{c} \in \{0,1\}^{\log s}} g(\mathbf{b}, \mathbf{c})$  where

$$g(\mathbf{b}, \mathbf{c}) = \widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\widetilde{W}_{i+1}(\mathbf{b}) + \widetilde{W}_{i+1}(\mathbf{c})) \\ + \widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\widetilde{W}_{i+1}(\mathbf{b}) \cdot \widetilde{W}_{i+1}(\mathbf{c}))$$

- A naïve implementation of **P** takes  $\Omega(S^3)$  time, where  $S$  is circuit size.
  - Same idea as prover implementation for #SAT protocol.
  - i.e., **P** evaluates  $g$  in each round of sum-check at all  $O(S^2/2^i)$  necessary points  $\mathbf{z}$ , taking  $O(S)$  time per point.

# GKR Prover Runtime: Details

Recall: Core of the GKR protocol is applying sum-check to compute  $\sum_{\mathbf{b}, \mathbf{c} \in \{0,1\}^{\log S}} g(\mathbf{b}, \mathbf{c})$  where

$$g(\mathbf{b}, \mathbf{c}) = \widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\widetilde{W}_{i+1}(\mathbf{b}) + \widetilde{W}_{i+1}(\mathbf{c})) \\ + \widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\widetilde{W}_{i+1}(\mathbf{b}) \cdot \widetilde{W}_{i+1}(\mathbf{c}))$$

- A naïve implementation of **P** takes  $\Omega(S^3)$  time, where  $S$  is circuit size.
  - Same idea as prover implementation for #SAT protocol.
  - i.e., **P** evaluates  $g$  in each round of sum-check at all  $O(S^2/2^i)$  necessary points  $\mathbf{z}$ , taking  $O(S)$  time per point.
- [CMT12]: **P** time is  $O(S \log S)$ .
  - Exploit structure in **multilinear** extensions  $\widetilde{\text{add}}_i$  and  $\widetilde{\text{mult}}_i$ . Ensures that each gate of  $C$  contributes to  $g(\mathbf{z})$  for  $O(1)$  relevant points  $\mathbf{z}$  in each round.



# GKR Prover Runtime: Details

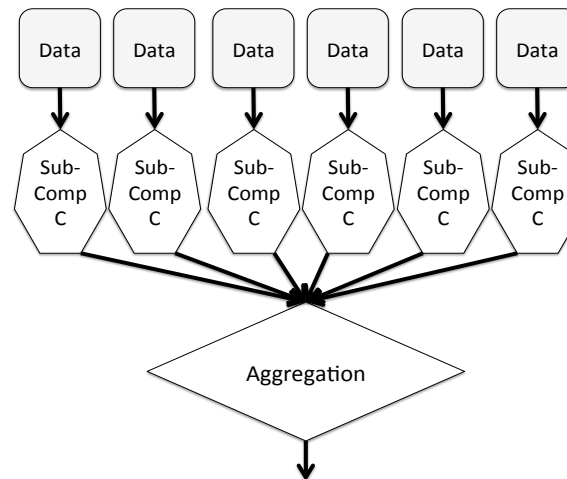
Recall: Core of the GKR protocol is applying sum-check to compute  $\sum_{\mathbf{b}, \mathbf{c} \in \{0,1\}^{\log S}} g(\mathbf{b}, \mathbf{c})$  where

$$g(\mathbf{b}, \mathbf{c}) = \widetilde{\text{add}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\widetilde{W}_{i+1}(\mathbf{b}) + \widetilde{W}_{i+1}(\mathbf{c})) \\ + \widetilde{\text{mult}}_i(\mathbf{r}_1, \mathbf{b}, \mathbf{c})(\widetilde{W}_{i+1}(\mathbf{b}) \cdot \widetilde{W}_{i+1}(\mathbf{c}))$$

- A naïve implementation of **P** takes  $\Omega(S^3)$  time, where  $S$  is circuit size.
  - Same idea as prover implementation for #SAT protocol.
  - i.e., **P** evaluates  $g$  in each round of sum-check at all  $O(S^2/2^i)$  necessary points  $\mathbf{z}$ , taking  $O(S)$  time per point.
- [CMT12]: **P** time is  $O(S \log S)$ .
  - Exploit structure in **multilinear** extensions  $\widetilde{\text{add}}_i$  and  $\widetilde{\text{mult}}_i$ . Ensures that each gate of  $C$  contributes to  $g(\mathbf{z})$  for  $O(1)$  relevant points  $\mathbf{z}$  in each round.
- All subsequent works seek to bring “Approach 3” to bear on the GKR protocol, letting **P** reuse work across rounds.

# GKR Prover Runtime: Details

- [Thaler13]:
  1. **P** time  $O(S)$  for circuits with “nice” wiring patterns.
  2. **P** time  $O(S \log S')$  for data parallel circuits  
i.e., that apply the same subcomputation (of size  $S'$ )  
independently to different pieces of data



# GKR Prover Runtime: Details

- [Thaler13]:
  1.  $P$  time  $O(S)$  for circuits with “nice” wiring patterns.
  2.  $P$  time  $O(S \log S')$  for data parallel circuits  
i.e., that apply the same subcomputation (of size  $S'$ )  
independently to different pieces of data
- [WJBSTWW17] improved the data parallel time to  $O(S + S' \log S')$ .

# GKR Prover Runtime: Details

- [Thaler13]:
  1.  $\mathbb{P}$  time  $O(S)$  for circuits with “nice” wiring patterns.
  2.  $\mathbb{P}$  time  $O(S \log S')$  for data parallel circuits  
i.e., that apply the same subcomputation (of size  $S'$ )  
independently to different pieces of data
- [WJBSTWW17] improved the data parallel time to  $O(S + S' \log S')$ .
- [ZGKPP18] extends to data parallel computations where each subcomputation may not be the same.

# GKR Prover Runtime: Details

- [Thaler13]:
  1.  $\mathbb{P}$  time  $O(S)$  for circuits with “nice” wiring patterns.
  2.  $\mathbb{P}$  time  $O(S \log S')$  for data parallel circuits  
i.e., that apply the same subcomputation (of size  $S'$ ) independently to different pieces of data
- [WJBSTWW17] improved the data parallel time to  $O(S + S' \log S')$ .
- [ZGKPP18] extends to data parallel computations where each subcomputation may not be the same.
- [XZZPS19] achieved  $O(S)$  time for general circuits.

# Rumination on Generality Vs. Efficiency

# Generality vs. Efficiency

- The GKR protocol for circuit evaluation has now been rendered optimally efficient for  $\mathbf{P}$  (up to constant factors).
- **Any** computation can be represented as a circuit evaluation (or satisfiability) problem.
  - But this can introduce tremendous overheads.
  - The GKR protocol **forces** the prover to compute the output in a prescribed manner, which may be far from optimal (gate-by-gate evaluation of a circuit).
- To achieve scalability, the gold standard is really **super-efficiency**.
  - i.e.,  $\mathbf{P}$  computed the right answer directly using the fastest known algorithm, and did a low-order amount of extra work to prove correctness