

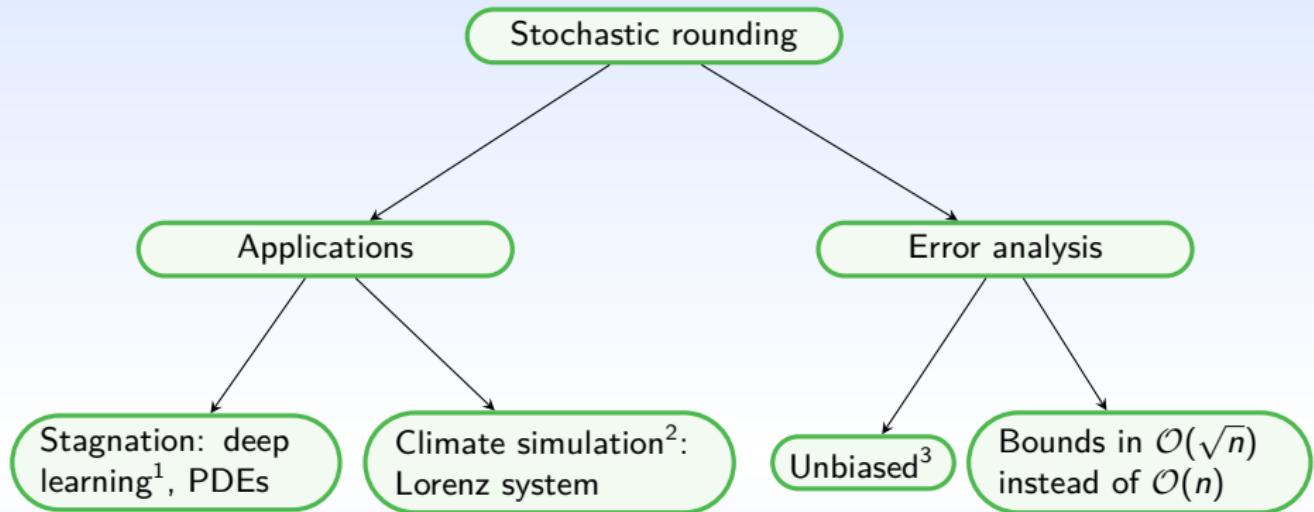
# Error Analysis for Algorithms using Stochastic Rounding

EL-Mehdi EL ARAR  
*el-mehdi.el-arar@inria.fr*

Journées InterFLOP/FPT-4



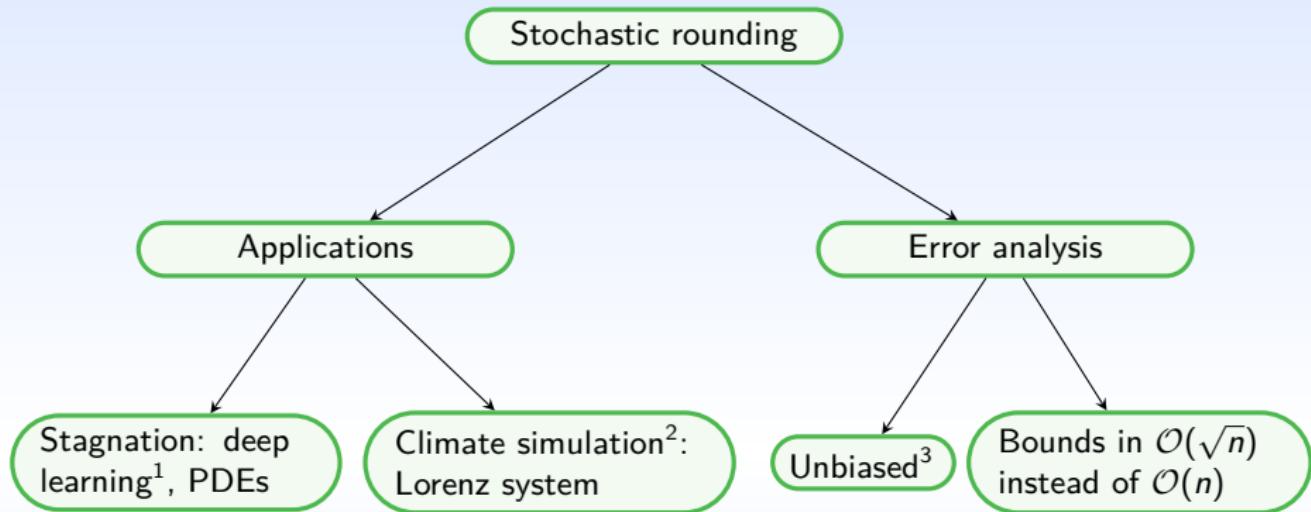
20/03/2025



<sup>1</sup>Gupta et al: Deep Learning with Limited Numerical Precision

<sup>2</sup>Paxton et al: Climate Modeling in Low Precision: Effects of Both Deterministic and Stochastic Rounding

<sup>3</sup>Parker: Monte Carlo Arithmetic: Exploiting Randomness in Floating-Point Arithmetic



**El arar:** Stochastic models for the evaluation of numerical errors

<sup>1</sup>Gupta et al: Deep Learning with Limited Numerical Precision

<sup>2</sup>Paxton et al: Climate Modeling in Low Precision: Effects of Both Deterministic and Stochastic Rounding

<sup>3</sup>Parker: Monte Carlo Arithmetic: Exploiting Randomness in Floating-Point Arithmetic

- ① Error Analysis of Algorithms with Stochastic Rounding
- ② Error Analysis for the Implementation of Stochastic Rounding in Hardware

# Stochastic rounding (SR)

- For  $x, y \in \mathbb{R}$  and  $\text{op} \in \{+, -, \times, /\}$

$$\text{SR}_p(x) = x(1 + \delta), \quad |\delta| \leq u_p$$

$$\text{SR}_p(x \text{ op } y) = (x \text{ op } y)(1 + \beta), \quad |\beta| \leq u_p$$

- Upward rounding  $\lceil x \rceil$  and downward rounding  $\lfloor x \rfloor$ :

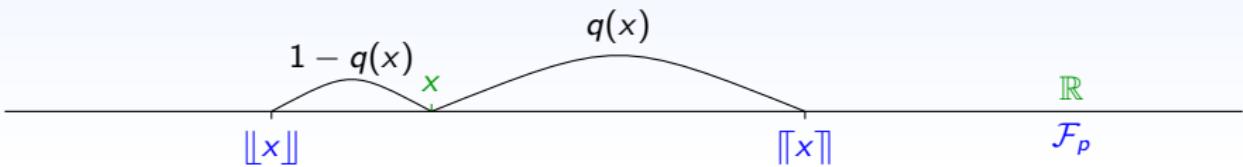


FIGURE.  $\text{SR}_p$  with  $q(x) = \frac{x - \lfloor x \rfloor}{\lceil x \rceil - \lfloor x \rfloor}$

- $\mathbb{E}(\text{SR}_p(x)) = q(x)\lceil x \rceil + (1 - q(x))\lfloor x \rfloor = x$ , then  $\mathbb{E}(\delta) = \mathbb{E}\left(\frac{\text{SR}_p(x) - x}{x}\right) = 0$

- Using  $SR_p$ , for  $x_1, x_2, x_3 \in \mathcal{F}_p$  and  $\text{op}_1, \text{op}_2 \in \{+, -, *, /\}$

$$c = x_1 \text{op}_1 x_2 \text{op}_2 x_3 \implies SR_p(c) = ((x_1 \text{op}_1 x_2)(1 + \delta_1) \text{op}_2 x_3)(1 + \delta_2),$$

with  $\mathbb{E}(\delta_1) = \mathbb{E}(\delta_2) = 0$

- Mean independence:  $\mathbb{E}[X_k/X_1, \dots, X_{k-1}] = \mathbb{E}(X_k)$  for all  $k$
- Independence  $\implies$  **Mean independence**  $\implies$  uncorrelatedness

## Lemma 1 (M. P. Connolly et al.).

For some  $\delta_1, \delta_2, \dots$ , *in that order* obtained from  $SR_p$ , the  $\delta_k$  are random variables with mean zero such that  $\mathbb{E}[\delta_k/\delta_1, \dots, \delta_{k-1}] = \mathbb{E}(\delta_k) = 0$ .

For  $z$  resulting of a multi-linear sum-product computation graph with  $n$  SR operations,

- $\Psi = \frac{\text{SR}_p(z) - z}{z}$  is a martingale
  - $E(\Psi) = 0$
  - $|\Psi|$  is bounded by  $\mathcal{O}(\sqrt{n}u_p)$  at fixed probability where  $n$  is the number of operations
- 
- The paper gives tighter bounds depending on the operations combinations.

Preprint: <https://arxiv.org/abs/2411.13601>

Error Analysis of sum-product algorithms under stochastic rounding  
de Oliveira Castro, El arar, Petit, Sohier

SR sum-product analysis gives error bounds for multi-linear algorithms:

- Inner product  $\mathcal{O}(\sqrt{n}u_p)$
- Horner's polynomial evaluation  $\mathcal{O}(\sqrt{n}u_p)$
- Pairwise summation  $\mathcal{O}(\sqrt{\log_2 n}u_p)$
- Karatsuba multiplication  $\mathcal{O}(\sqrt{\log_2 n}u_p)$

## Example: inner product

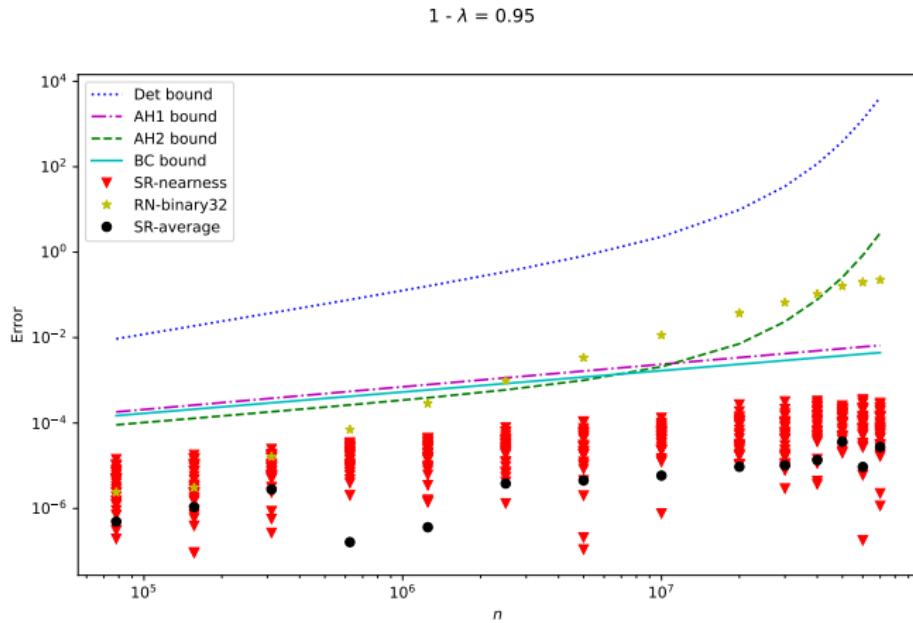


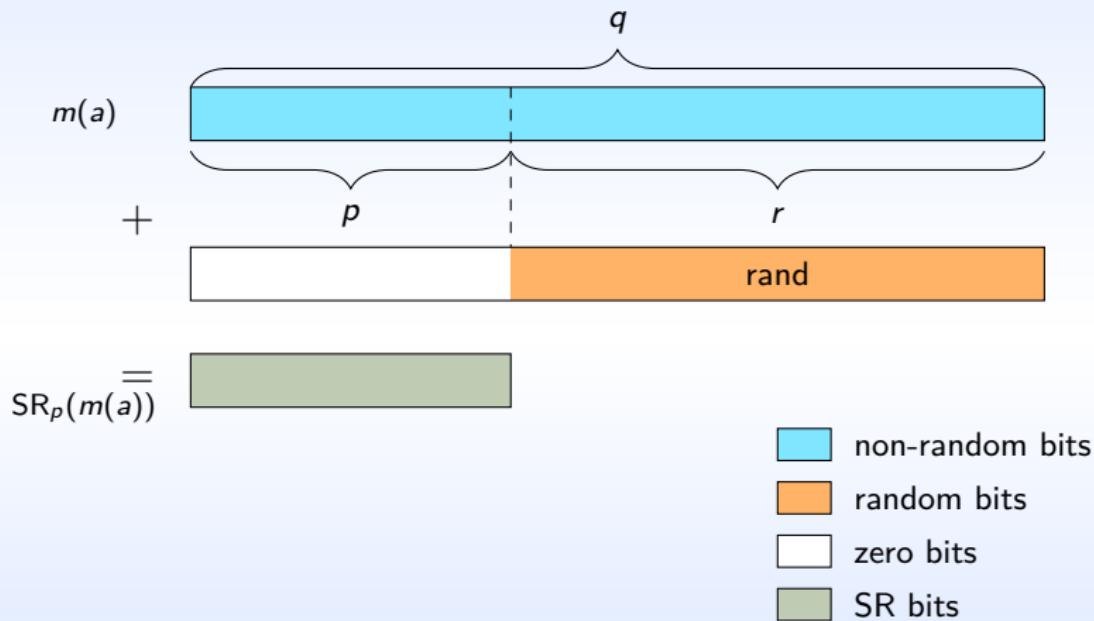
FIGURE. Probabilistic bounds with probability  $1 - \lambda = 0.95$  vs deterministic bound of the computed forward errors of the inner product for floating-point numbers chosen uniformly at random in  $[0; 1]$  with  $u_p = 2^{-23}$

- 1 Error Analysis of Algorithms with Stochastic Rounding
- 2 Error Analysis for the Implementation of Stochastic Rounding in Hardware

# How can we implement this in hardware?

## Low precision

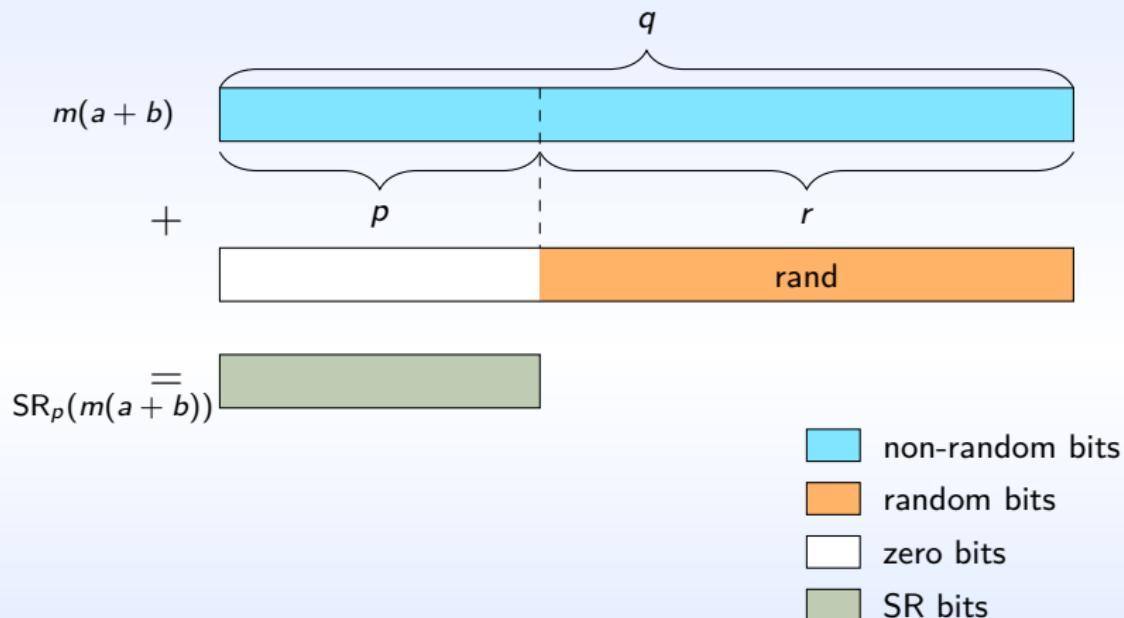
Let  $a \in \mathcal{F}_q$ , if we represent  $a$  in  $\mathcal{F}_p$  such that  $p < q$ , we take  $r = q - p$ .



# How can we implement this in hardware?

## Addition

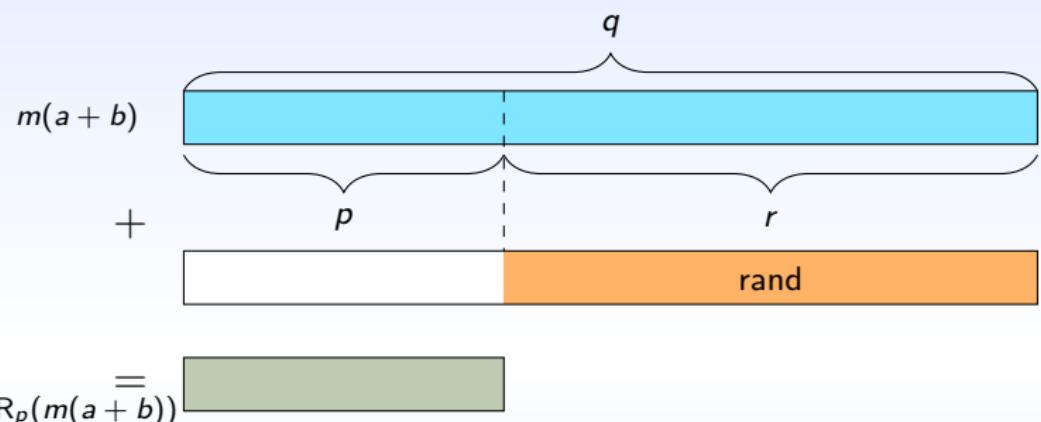
Let  $a, b \in \mathcal{F}_p$ , if we compute  $a + b$  in  $\mathcal{F}_q$  such that  $p < q$ , we take  $r = q - p$ .



# How can we implement this in hardware?

## Addition

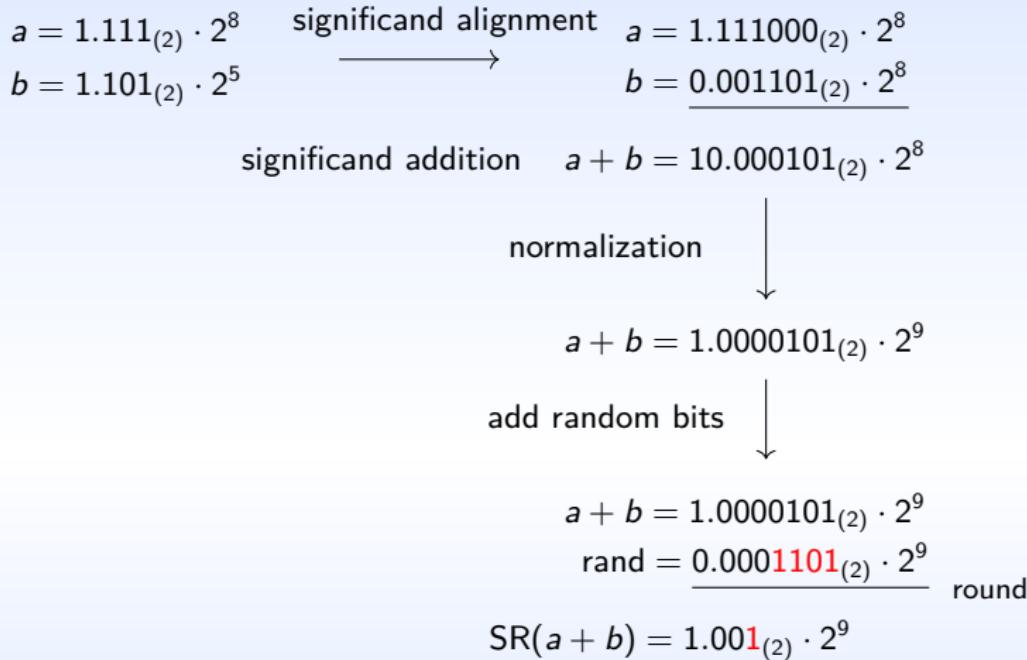
Let  $a, b \in \mathcal{F}_p$ , if we compute  $a + b$  in  $\mathcal{F}_q$  such that  $p < q$ , we take  $r = q - p$ .



A Stochastic Rounding-Enabled Low-Precision Floating-Point MAC for DNN Training  
"Ali, Sami Ben and Filip, Silviu-Ioan and Sentieys, Olivier"

- non-random bits
- random bits
- zero bits
- SR bits

## Toy example: $a + b$



## Toy example: $a + b$

$$\begin{array}{l} a = 1.111_{(2)} \cdot 2^8 \\ b = 1.101_{(2)} \cdot 2^5 \end{array} \xrightarrow{\text{significand alignment}} \begin{array}{l} a = 1.111000_{(2)} \cdot 2^8 \\ b = \underline{0.001101}_{(2)} \cdot 2^8 \end{array}$$

significand addition     $a + b = 10.000101_{(2)} \cdot 2^8$

normalization



$$a + b = 1.0000101_{(2)} \cdot 2^9$$

add random bits

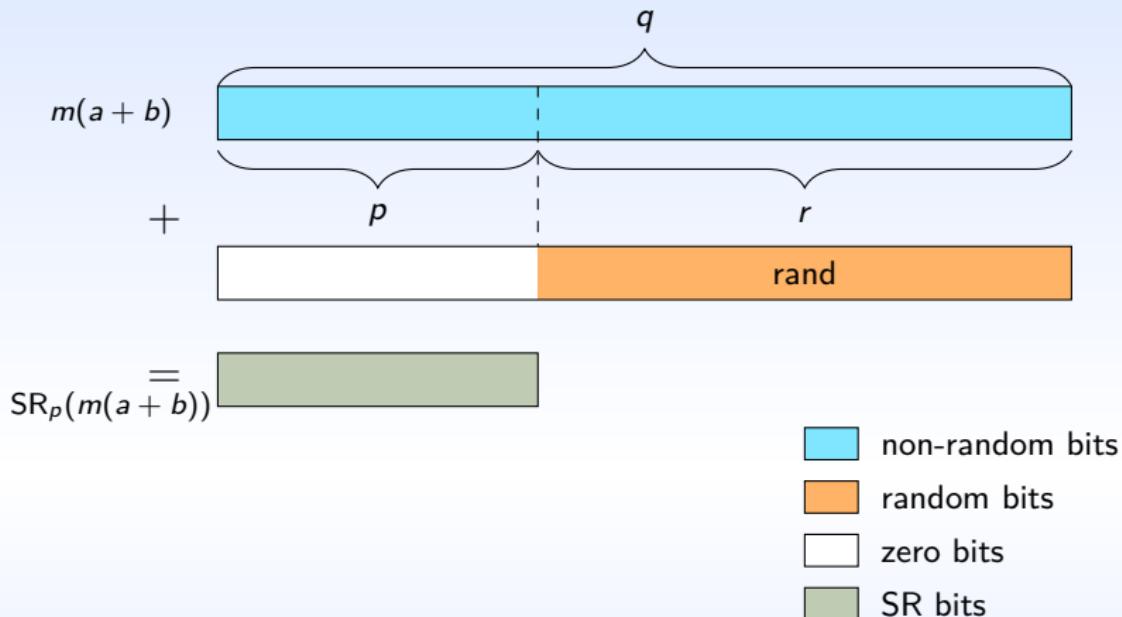


$$a + b = 1.0000101_{(2)} \cdot 2^9$$

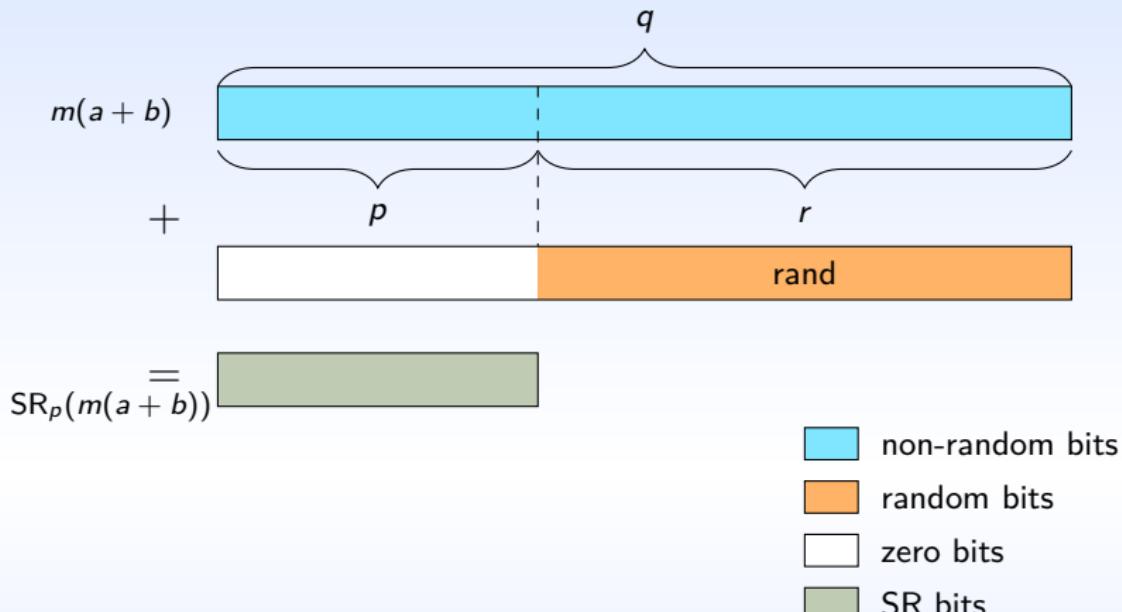
$$\begin{array}{r} \text{rand} = \underline{0.000\textcolor{red}{0110}}_{(2)} \cdot 2^9 \\ \text{round} \end{array}$$

$$\text{SR}(a + b) = 1.000_{(2)} \cdot 2^9$$

## How can we implement this in hardware?

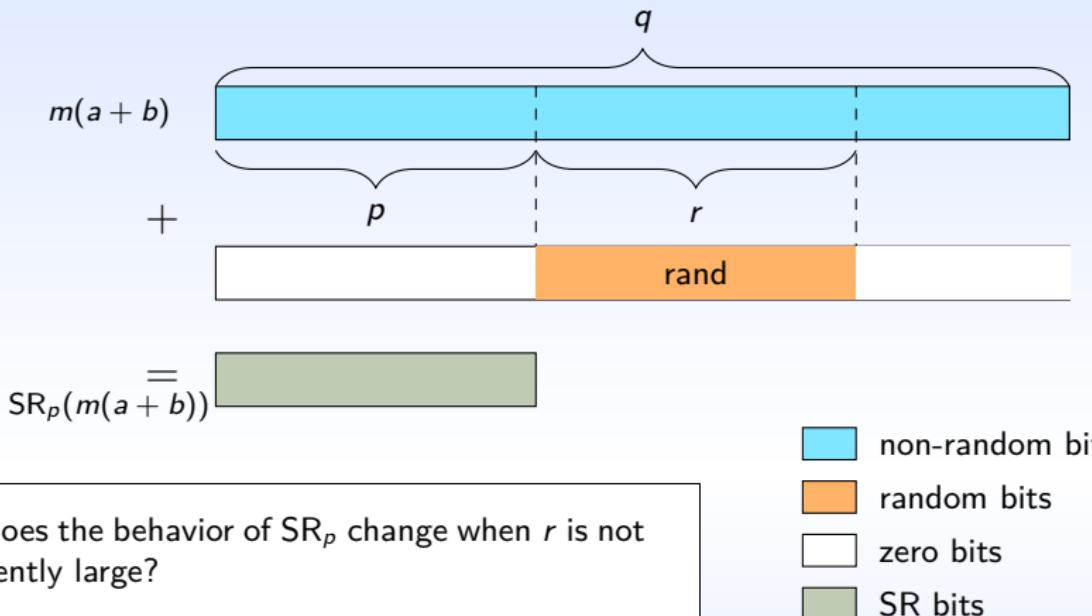


## How can we implement this in hardware?



Expensive!!

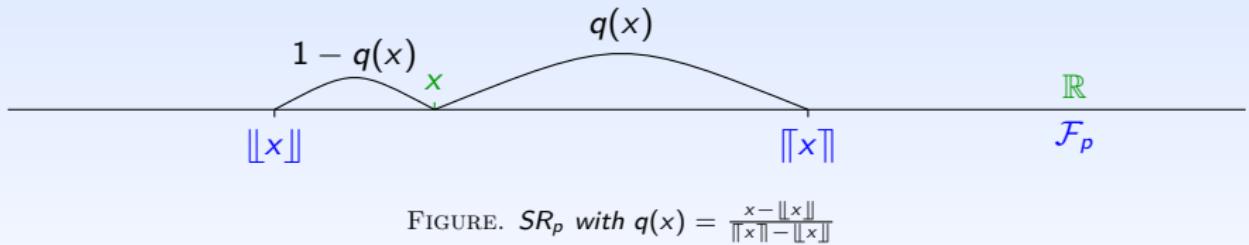
## How can we implement this in hardware?



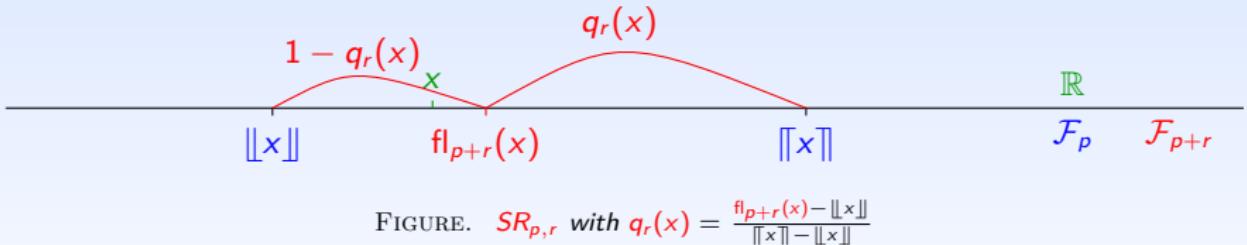
How does the behavior of  $\text{SR}_p$  change when  $r$  is not sufficiently large?

Expensive!!

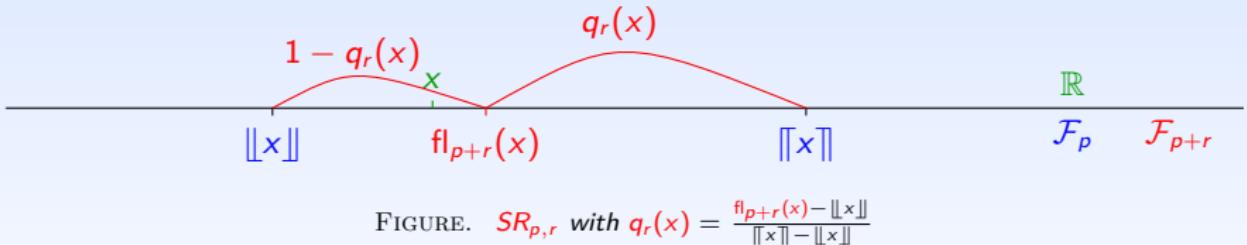
# Limited-precision stochastic rounding



# Limited-precision stochastic rounding



# Limited-precision stochastic rounding



- $\text{SR}_{p,r}(x) = x(1 + \delta)$ ,  $|\delta| \leq u_p \neq \text{fl}_{p+r}(x) = x(1 + \beta)$ ,  $|\beta| \leq u_{p+r}$
- $\mathbb{E}(\text{SR}_{p,r}(x)) = q_r(x)\lceil x \rceil + (1 - q_r(x))\lfloor x \rfloor = \text{fl}_{p+r}(x)$ , then

$$\mathbb{E}(\delta) = \mathbb{E}\left(\frac{\text{SR}_{p,r}(x) - x}{x}\right) = \beta$$

- The mean independence is lost

$$\mathbb{E}(\delta_k | \delta_1, \dots, \delta_{k-1}) = \beta_k \neq \mathbb{E}(\delta_k)$$

- $\beta_k$  is a random variable and  $\mathbb{E}(\beta_k) = \mathbb{E}(\delta_k)$

## Doob–Meyer decomposition

### Lemma 2.

Let  $\delta_1, \delta_2, \dots, \delta_n$  be random errors produced by a sequence of elementary operations using  $\text{SR}_{p,r}$ , and let  $\beta_1, \beta_2, \dots, \beta_n$  be their corresponding errors incurred by  $\text{fl}_{p+r}$ . Then, the random variables  $\alpha_k = \delta_k - \beta_k$  for  $1 \leq k \leq n$  are mean independent

$$\mathbb{E}(\alpha_k \mid \alpha_1, \dots, \alpha_{k-1}) = \mathbb{E}(\alpha_k) = 0.$$

Moreover, for all  $1 \leq i \leq n$

$$\prod_{k=i}^n (1 + \delta_k) = \prod_{k=i}^n (1 + \alpha_k) + \mathcal{B}_i,$$

with

$$|\mathcal{B}_i| \leq \gamma_{n-i+1}(u_p + u_{p+r}) - \gamma_{n-i+1}(u_p)$$

## Theorem 3.

For  $y = \sum_{i=1}^n a_i b_i$  and  $0 < \lambda < 1$ , the quantity  $SR_{p,r}(y)$  satisfies

$$\begin{aligned}\frac{|SR_{p,r}(y) - y|}{|y|} &\leq \kappa(a \circ b) \left( \sqrt{u_p \gamma_{2n}(u_p)} \sqrt{\ln(2/\lambda)} + \gamma_n(u_p + u_{p+r}) - \gamma_n(u_p) \right) \\ &= \kappa(a \circ b) \left( \sqrt{2n} \sqrt{\ln(2/\lambda)} u_p + \textcolor{red}{nu_{p+r}} \right) + \mathcal{O}(\|(u_p, u_{p+r})\|_2)\end{aligned}$$

with probability at least  $1 - \lambda$ .

- It can be applied to all previous algorithms studied with  $SR_p$

## Error analysis of algorithms with limited-precision SR

$$\prod_{k=i}^n (1 + \delta_k) = \prod_{k=i}^n (1 + \alpha_k) + \mathcal{B}_i$$

Martingale

$$|\mathcal{B}_i| \leq \gamma_{n-i+1}(u_p + u_{p+r}) - \gamma_{n-i+1}(u_p)$$

$$\mathcal{O}(\sqrt{n}u_p)$$

$$\mathcal{O}(nu_{p+r})$$

$$\mathcal{O}(\sqrt{n}u_p + nu_{p+r})$$

- We have the same result with BC method

## Lemma 4.

*Theorem 3 leads to*

$$\mathcal{O}(\sqrt{n}u_p + nu_{p+r})$$

*we want*

$$\sqrt{n}u_p > nu_{p+r}$$

*we thus have this good rule of thumb*

$$r \geq \lceil (\log_2 n)/2 \rceil$$

- Rosenbrock function (**srfloat**)
  - Parameter update in deep neural network training (**MPTorch**)
- 
- <https://github.com/MehdiElArar/srfloat>
  - <https://github.com/mptorch/mptorch>

## Numerical experiments: Rosenbrock function

The Rosenbrock function is a non-convex function defined by

$$f(x_1, x_2) = (1 - x_1)^2 + 100(x_2 - x_1^2)^2$$

with a global minimum of 0, occurring at  $\mathbf{x}^* = (1, 1)$ .

The gradient descent:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)$$

# Rosenbrock function

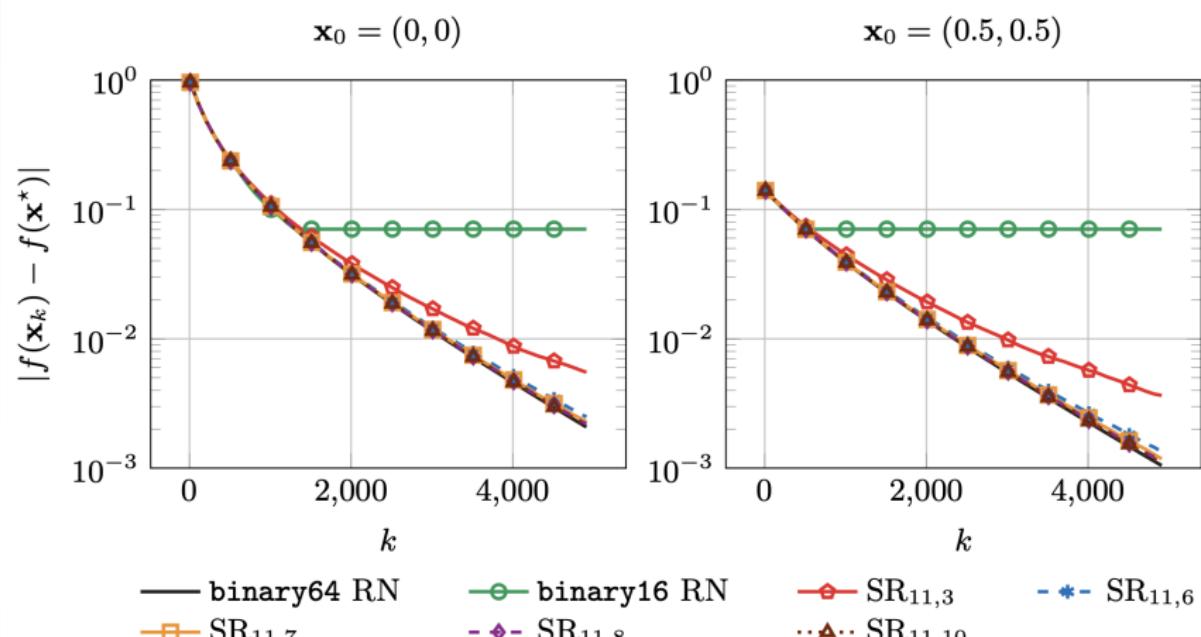


FIGURE. Convergence profiles for 5,000 iterations of gradient descent on the Rosenbrock function. For both experiments, we average each  $SR_{11,r}$  error over 500 different runs, and the learning rate is  $t_k = 0.001$ .

# Rosenbrock function

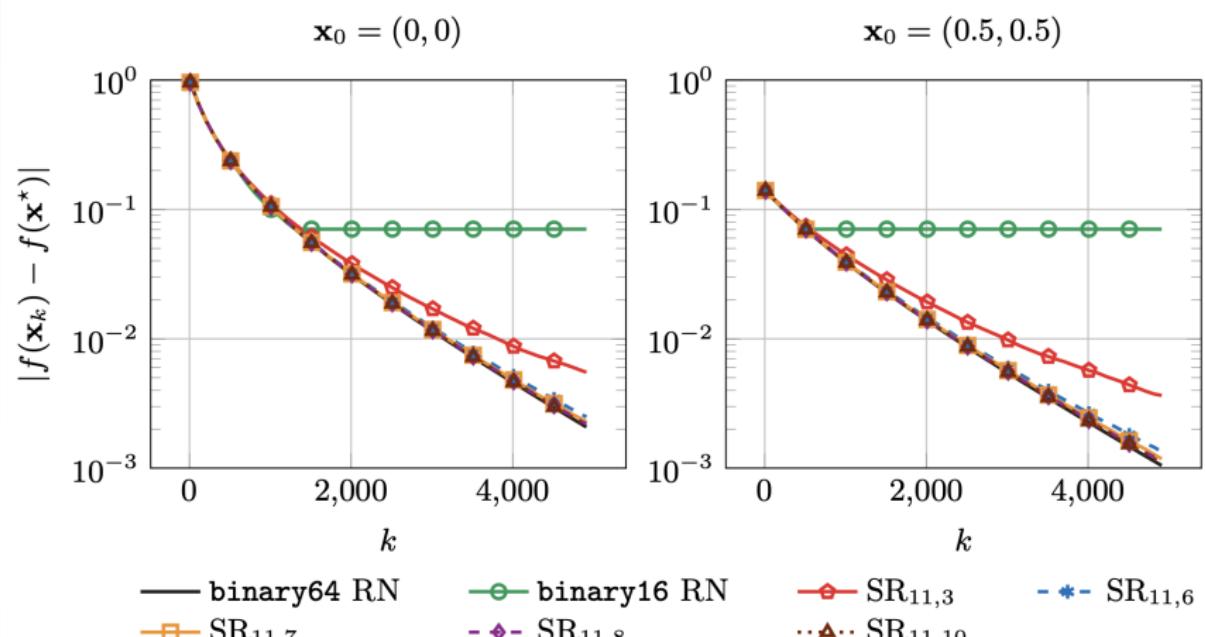


FIGURE. Convergence profiles for 5,000 iterations of gradient descent on the Rosenbrock function. For both experiments, we average each  $SR_{11,r}$  error over 500 different runs, and the learning rate is  $t_k = 0.001$ .

$$\lceil \log_2(5,000)/2 \rceil = 7$$

## Focus:

Training a ResNet32<sup>1</sup> model on the CIFAR-10 dataset

## Training Setup:

- Hyperparameters:

- ▶ Batch Size: 128, Momentum:  $\mu = 0.9$
- ▶ Total Training: 64,000 iterations (200 epochs)
- ▶ Learning Rate:  $t_k = 0.1$ , reduced by 10 at 32,000 and 48,000 iterations

## Numerical Precision:

- Arithmetic: bfloat16 ( $p = 8$ )

- Update Rule:

$$\begin{aligned}\mathbf{v}_{k+1} &= \circ(\mu \mathbf{v}_k + \mathbf{g}_k), \\ \mathbf{x}_{k+1} &= \circ(\mathbf{x}_k - t_k \mathbf{v}_{k+1})\end{aligned}$$

- Components:

- ▶  $\mathbf{v}_k$ : Velocity vector
- ▶  $\mathbf{g}_k$ : Gradient of the loss function

---

<sup>1</sup>Deep Residual Learning for Image Recognition

# Parameter update in deep neural network training

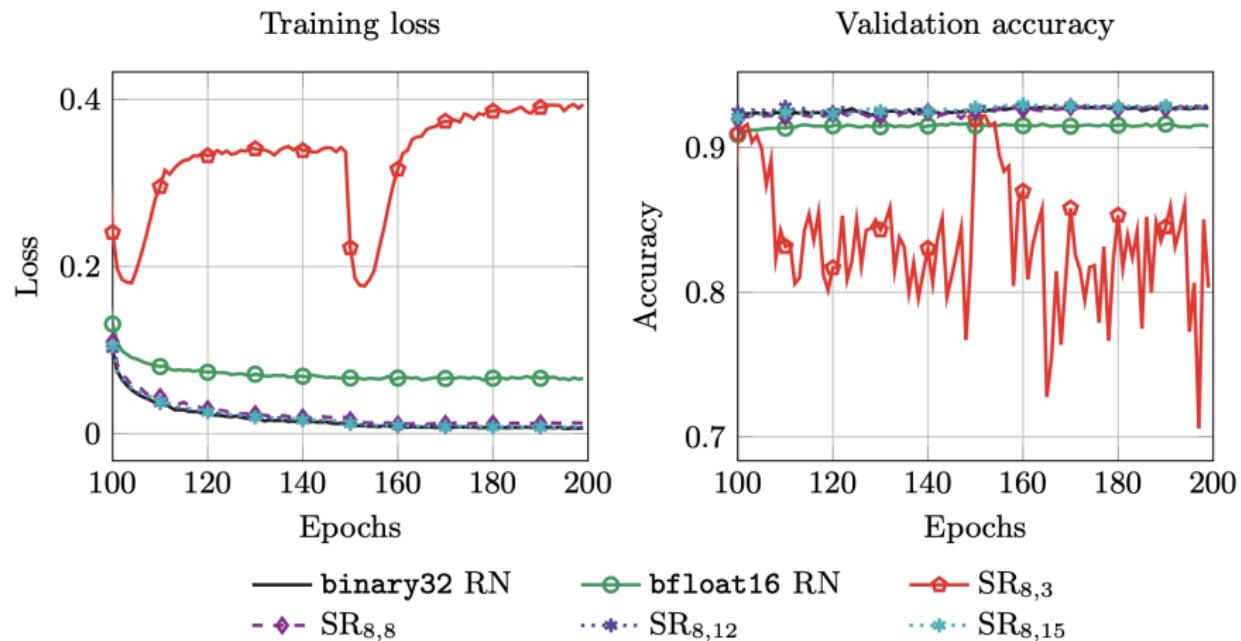


FIGURE. In the baseline configuration, binary32 arithmetic with RN is used for computing, and the same format is used for storage. For the low-precision configurations, parameters are stored and updated using bfloat16 arithmetic with either RN or  $SR_{p,r}$ .

# Parameter update in deep neural network training

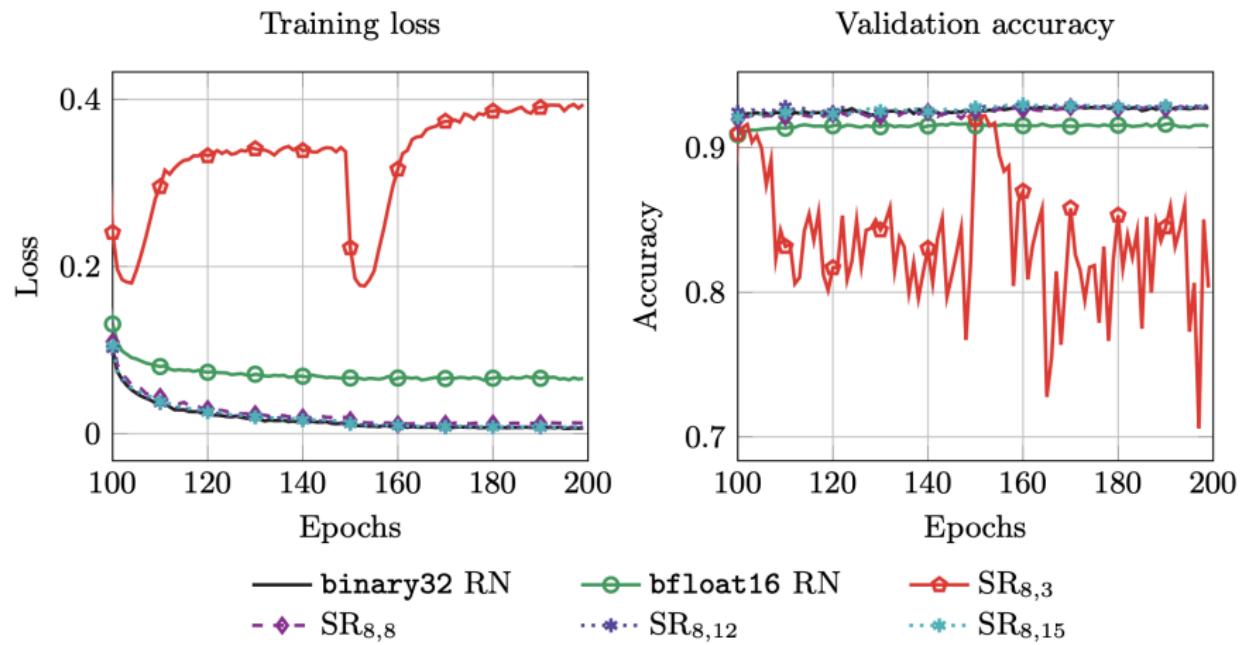


FIGURE. In the baseline configuration, *binary32* arithmetic with *RN* is used for computing, and the same format is used for storage. For the low-precision configurations, parameters are stored and updated using *bfloat16* arithmetic with either *RN* or  $SR_{p,r}$ .

$$\lceil \log_2(64,000)/2 \rceil = 8$$

# Summary

	SR <sub>p</sub>	SR <sub>p,r</sub>
Unbiased	✓	✗
Mean independence	✓	✗
Probabilistic bound	$\mathcal{O}(\sqrt{n}u_p)$	$\mathcal{O}(\sqrt{n}u_p + nu_{p+r})$
Rule of thumb		$r \geq \lceil (\log_2 n)/2 \rceil$

TABLE. *Classic stochastic rounding versus limited-precision stochastic rounding*

Preprint: <https://arxiv.org/abs/2408.03069>

Probabilistic error analysis of limited-precision stochastic rounding  
**El-Mehdi El arar**, Massimiliano Fasi, Silviu-Ioan Filip, Mantas Mikaitis