# Custom Tailored FPGA Boson Sampling

Gregory Morse [1]
morse@inf.elte.hu

Tamás Kozsik [1]
kto@elte.hu

Péter Rakyta [2]
peter.rakyta@ttk.
elte.hu

[1] Department of Programming Languages and Compilers,

[2] Department of Physics of Complex Systems, Eötvös Loránd
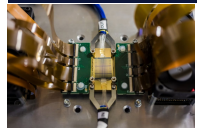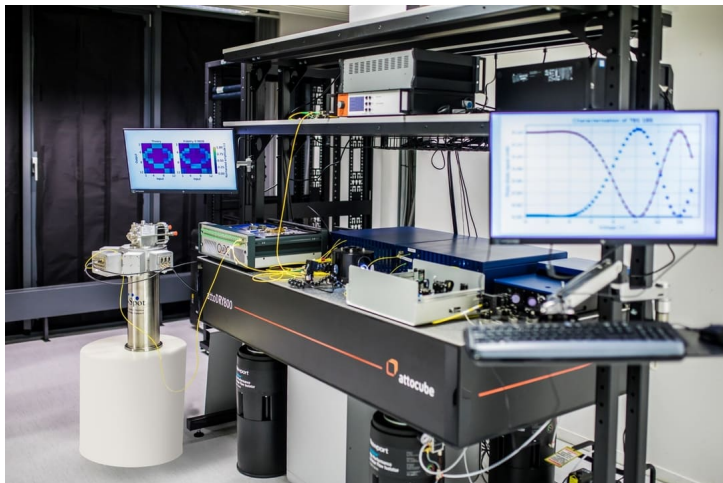Tudományegyetem/University (ELTE), Budapest, Hungary

# Introduction

- Computing the permanent of a matrix finds an important application in the context of boson sampling
- BB/FG permanent formula with a reflected binary Gray code $\mathcal{O}(n.2^{n-1})$
- Run it in parallel on 4 SLRs (Super Logic Regions) $\mathcal{O}(n.2^{n-3})$
- Up to 40x40 matrix permanents @ 280MHz
- Dual FPGA - specialized kernel with twice as fast operation
- Generalize to repeated rows/columns up to 40 photons @ 240MHz

Keywords: Boson Sampling, Matrix Permanent, FPGA, dataflow, Repeated Row and Column Permanent, Reflected N-ary Gray code

# Real World Boson Sampling Setup



- QuiX Photonic Quantum Computer

# Optimally Efficient Classical Permanent Algorithms

- BB/FG formula:

$$\text{perm}(A) = \frac{1}{2^{m-1}} \sum_{\delta} \left( \prod_{k=1}^{m} \delta_k \right) \prod_{j=1}^{m} \sum_{i=1}^{m} \delta_i a_{i,j}, \tag{1}$$

where $A$ is an $m \times m$ square matrix describing the interferometer and $\delta$ is a binary Gray code. Adaptable to sub-computations of repeated-row rectangular permanents.

- With independent repeated rows and columns:

$$\text{perm}(A, M, N) =$$

$$\frac{1}{2^{n-1}} \sum_{\Delta} \left( \prod_{k=1}^{m} (-1)^{\Delta_k} \binom{M_k}{\Delta_k} \right) \prod_{j=1}^{m} \left( \sum_{k=1}^{m} (M_k - 2\Delta_k) a_{k,j} \right)^{N_j} \tag{2}$$
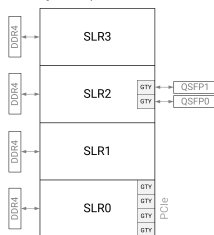
$M$ and $N$ are the row and column multiplicities respectively such that the photon count $n = \sum M_i = \sum N_j$ and $\Delta$ is the n-ary Gray code, required for efficient computation.

PCIExpress clock operates at 250MHz, initialization limitation 16nm
Ultrascale+ architecture, Vivado compiler supporting up to 500MHz



Figure 3: **Floorplan of the XCU250 Device**

| Specification | U250 | |
|---|---|---|
| | **Active Cooling Version** | **Passive Cooling Version** |
| Product SKU | A-U250-A64G-PQ-G | A-U250-P64G-PQ-G |
| Thermal cooling solution | Active | Passive |
| Weight | 1122g | 1066g |
| Form factor | Full height, full length, dual width | Full height, ¾ length, dual width |
| Total electrical card load† | 215W | |
| Network interface | 2x QSFP28 | |
| PCIe Interface | Gen3 x16 | |
| Look-up tables (LUTs) | 1,728K | |
| Registers | 3,456K | |
| DSP slices | 12,288 | |
| UltraRAMs | 1,280 | |
| DDR total capacity | 64 GB | |
| DDR maximum data rate | 2400 MT/s | |
| DDR total bandwidth | 77 GB/s | |

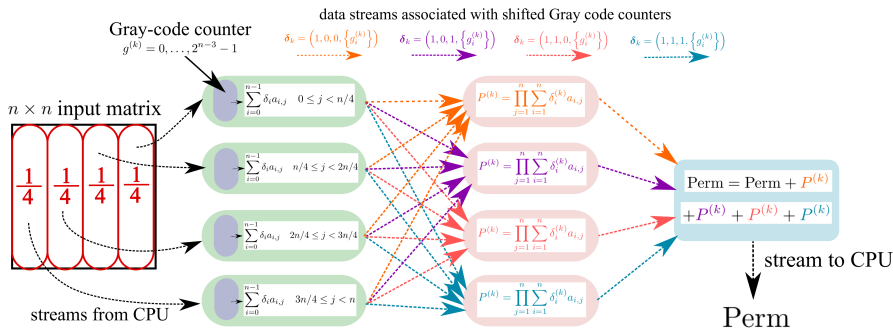# Accuracy of single and double precision floating vs. normalized fixed point

- We developed a GNU MPFR (multi-precision floating point reliability library) infinite precision wrapper to check accuracy with realistic data as part of **Piqasso Boost** extension to piquasso PIQUASSO

- Single/double precision CPU variants use 4M+2A complex multiplication $(a + bi)(c + di) = (ac - bd) + (bc + ad)i$

- Infinite Precision and FPGA variants use Knuth 3M+5A complex multiplication $x = c(a + b), (x - b(c + d)) + (x + a(d - c))i$

- Complex number normalization computed by worst-case column sums using a Euclidean vector inspired technique, keeping all computations $-1 \leq a, b, |a + bi| \leq 1$

- Outer sum of BB/FG precise number of bits determined
$$\frac{\max(\sum_{k=0}^{\lfloor (n-1)/2 \rfloor} (2(2k+2)-n)^n \binom{n-1}{2k+1}, \sum_{k=0}^{\lfloor (n-1)/2 \rfloor} (2(2k+1)-n)^n \binom{n-1}{2k})}{n^n}$$ for $n = 40$, then $\pm 27$ is maximum partial sum requiring 6 integer bits (including sign)

# Design

- Time complexity: $\mathcal{O}(2^{n-1-k})$ where $k = 2$ for single, $k = 3$ for dual
- Area: $\mathcal{O}(n^2)$ dominated by matrix storage in FFs (flip-flops)
- Multiplication area based on a product tree with Karatsuba rectangular tiling to match the 18x25 signed DSP multipliers of the FPGA $\mathcal{O}(b \log b)$ where tree depth is 6: 20 (b=64-bit) -> 10 (93-bit) -> 5 (110-bit) -> 2 (127-bit) -> 1 (127-bit) -> 1 (127-bit)

# Repeated Row/Column Design

- Reflected N-ary Gray code (using direction encoding (DE))
- Binomial coefficients computed with a loop length=9, division by "magic number" multiplication
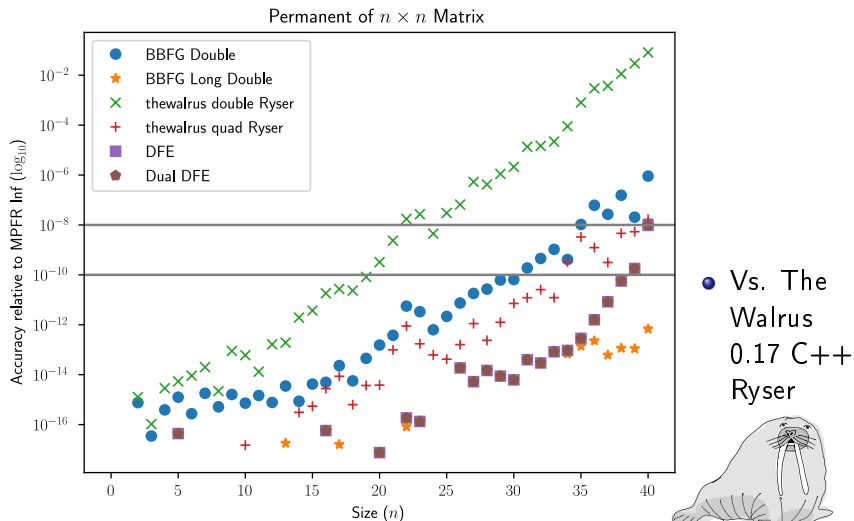  - Incremental update by Gray code decreasing $b_i = \frac{b_{i-1} \times k}{n-k+1}$ otherwise $b_i = \frac{b_{i-1} \times (n-k)}{k+1}$
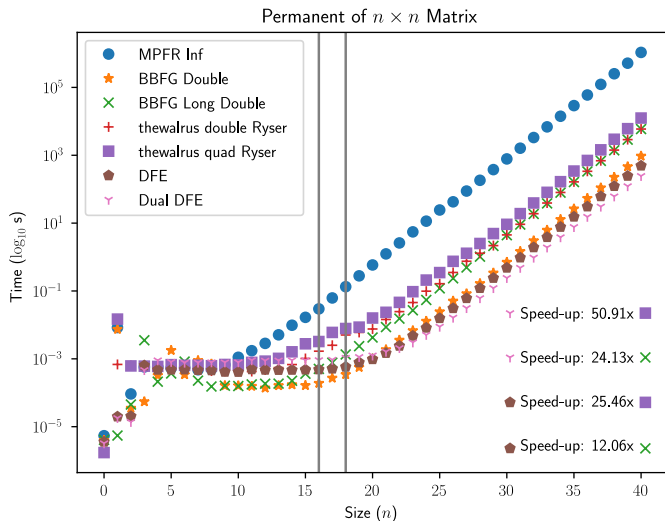- Staggering the Gray code at even intervals

| Counter Chain | DEGC | Gray Code (GC) | $b_i$ |
|:---:|:---:|:---:|:---:|
| **(0, 0, 0)** | **(0, 0, 0)** | **(0, 0, 0)** | **1** |
| (1, 0, 0) | (1, 0, 0) | (1, 0, 0) | 1 |
| **(0, 1, 0)** | **(2, 1, 0)** | **(1, 1, 0)** | **2** |
| (1, 1, 0) | (3, 1, 0) | (0, 1, 0) | 2 |
| **(0, 2, 0)** | **(0, 2, 0)** | **(0, 2, 0)** | **1** |
| (1, 2, 0) | (1, 2, 0) | (1, 2, 0) | 1 |
| **(0, 0, 1)** | **(2, 3, 1)** | **(1, 2, 1)** | **2** |
| (1, 0, 1) | (3, 3, 1) | (0, 2, 1) | 2 |
| **(0, 1, 1)** | **(0, 4, 1)** | **(0, 1, 1)** | **4** |
| (1, 1, 1) | (1, 4, 1) | (1, 1, 1) | 4 |
| **(0, 2, 1)** | **(2, 5, 1)** | **(1, 0, 1)** | **2** |
| (1, 2, 1) | (3, 5, 1) | (0, 0, 1) | 2 |
| **(0, 0, 2)** | **(0, 0, 2)** | **(0, 0, 2)** | **1** |
| (1, 0, 2) | (1, 0, 2) | (1, 0, 2) | 1 |
| **(0, 1, 2)** | **(2, 1, 2)** | **(1, 1, 2)** | **2** |
| (1, 1, 2) | (3, 1, 2) | (0, 1, 2) | 2 |
| **(0, 2, 2)** | **(0, 2, 2)** | **(0, 2, 2)** | **1** |
| (1, 2, 2) | (1, 2, 2) | (1, 2, 2) | 1 |

Table: Example for Non-Anchor Row Multiplicities (1, 2, 2) with 3x3x2=18 values
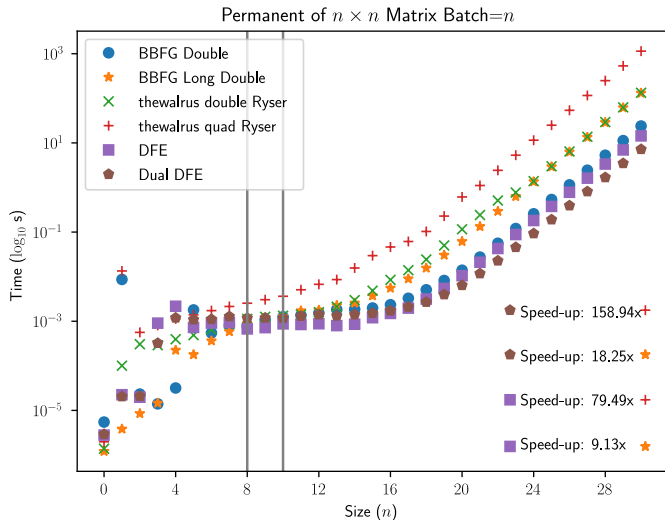
Permanent of $n \times n$ Matrix

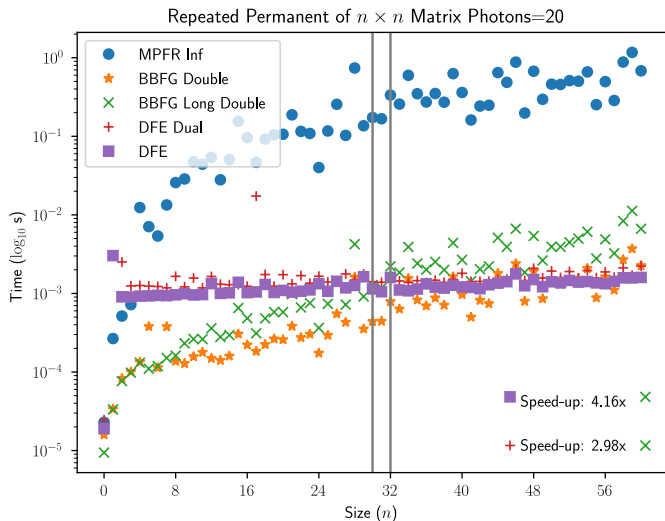- Vs. The Walrus 0.17 C++ Ryser

Permanent of $n \times n$ Matrix

- Initialization delay crossover threshold for single and dual marked based on precise long double calculators

# FPGA Batching Advantage



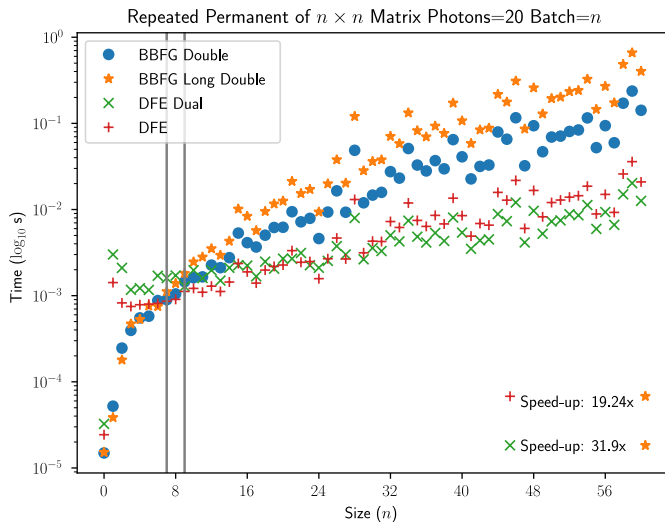Permanent of $n \times n$ Matrix Batch=$n$

- Batches reduce the cross-over threshold
- Kernels designed for automatic control signal resetting/-counter wrapping

# Performance of 240MHz for repeated permanents



Repeated Permanent of $n \times n$ Matrix Photons=20

- MPFR Inf
- BBFG Double
- BBFG Long Double
- DFE Dual
- DFE

Time ($\log_{10}$ s)
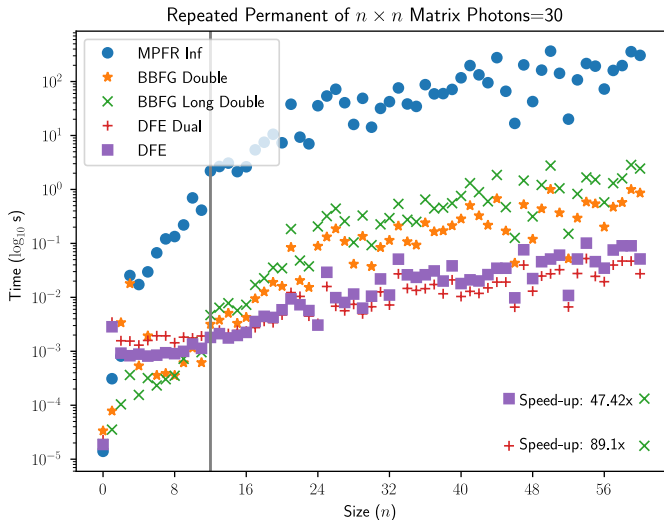
Size ($n$)

Speed-up: 4.16x
Speed-up: 2.98x

- Significant initialization delay and high crossover threshold
- At 30 photons, this no longer occurs

# FPGA Batching Advantage for repeated permanents



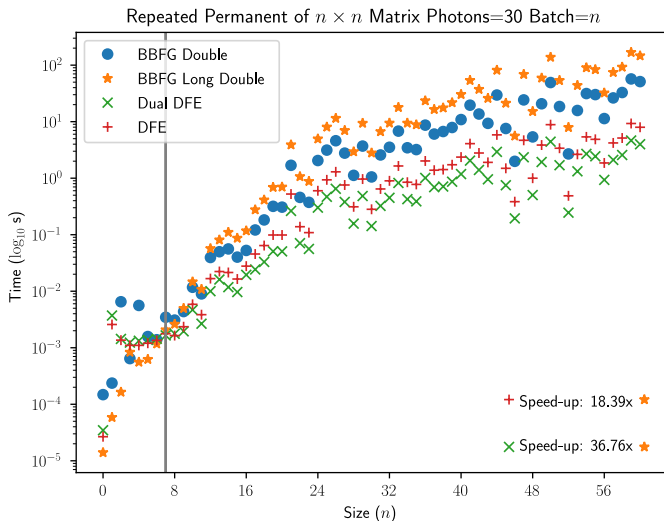Repeated Permanent of $n \times n$ Matrix Photons=20 Batch=$n$

- Batch size $n$ based on realistic use case Faster classical Boson Sampling (Clifford and Clifford 2018)
- Ouput state/row multiplicity fixed over batch of changing input states/column multipicities

# Performance of 240MHz for repeated permanents



Repeated Permanent of $n \times n$ Matrix Photons=30

Legend:
- MPFR Inf
- BBFG Double
- BBFG Long Double
- DFE Dual
- DFE

Speed-up: 47.42x
Speed-up: 89.1x

Axis: Time ($\log_{10}$ s) vs Size ($n$)

- DFE variants have clearly become best performing

# FPGA Batching Advantage for repeated permanents



Repeated Permanent of $n \times n$ Matrix Photons=30 Batch=$n$

- BBFG Double
- BBFG Long Double
- Dual DFE
- DFE

Speed-up: 18.39x

Speed-up: 36.76x

- The benefit is further demonstrated

# Time, Area and Power Analysis

- FPGA image upload time: 56.2 seconds for single, 112.9 seconds for dual
- Actual runtime: $t = t_0 + \frac{n-1+2^{n-1-k}}{f}$ so for 280MHz dual on a 40x40 matrix: $\frac{40-1+2^{36}}{280\times1000000} = 245$ seconds
- Effective equivalent: $\frac{(C_A+C_M)*280*10^6}{10^9}$ where $C_A = 2A = 2*(40+4)$ and $C_M = 4M + 2A = 6*4*39$ represent complex addition and multiplication respectively, yielding 285.5 GFLOPS for single mode
- Power estimate (in KWh): $\frac{w*t}{60*60*1000}$ so for 280MHz single on a 40x40 matrix: $\frac{14.83\times490}{3600000} = 0.002$KWh

```
FINAL POWER REPORT
Total On-Chip Power (W) 14.83 (budget: 135.00)
Dynamic Power (W)         11.56
Device Static Power(W)     3.27
FINAL RESOURCE USAGE
FPGA: xcU250-FIGD2104-2L-E
Logic utilization:      2029689 / 5184000 (39.15%)
  LUTs:                  741862 / 1728000 (42.93%)
  Primary FFs:          1287827 / 3456000 (37.26%)
DSP blocks:                8304 /   12288 (67.58%)
Block memory (BRAM18):      647 /    5376 (12.03%)
Block memory (URAM):        126 /    1280 ( 9.84%)
```
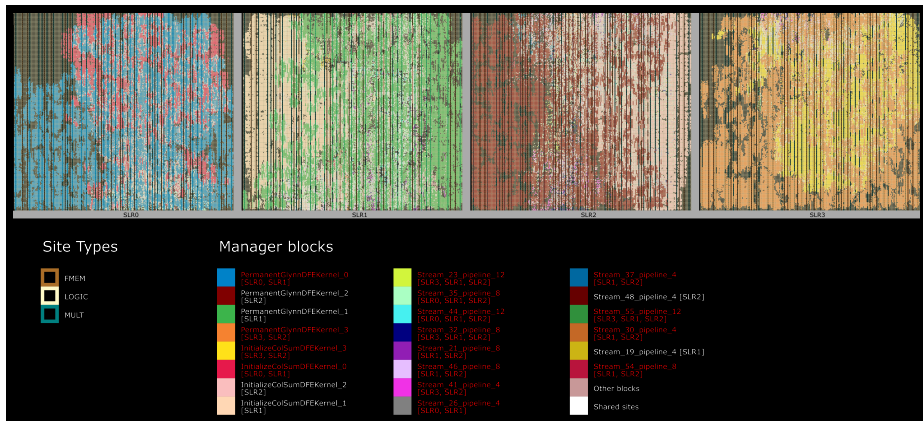
# Aerial View



- After careful pipelining and reducing fanout, streams become the largest routing issue

# Conclusion and Future Research

- FPGAs are competitive for multiplication intensive implementations
- State-of-the-art algorithms are often needed for maximizing resources, cannot rely on default implementations
- Computation of the Loop Hafnian using GroqCard in concert with data preparation on the FPGA



Active Option