

Schrodinger Equation Accelerator on PYNQ via HLS



On board test by PYNQ-Z2

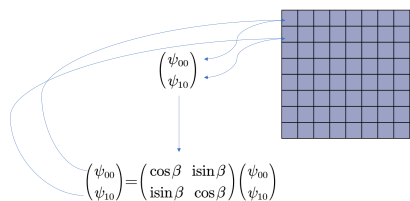
INTRODUCTION

Traditional numerical approaches for solving a PDE, such as the Runge-Kutta methods, are not suitable for quantum settings because they do not **numerically preserve the probability** (sum of the squared modulus of the wave function):

$$\int |\varphi|^2 dr$$

quantum mechanics should always be one.

To preserve the probability, we should turn to other algorithms, such as the Chebyshev algorithm and the Trotter-Suzuki algorithm. The accelerator is based on the **Trotter-Suzuki algorithm** [1], in which the transformation matrix from $t=0$ to $t=\Delta t$ has a determinant of one, meaning that it is a proper rotation matrix and will not alter the squared modulus of the vector it is acting on.



Principle of Trotter-Suzuki algorithm

[1] De Raedt, Hans. "Computer simulation of quantum phenomena in nanoscale devices." *Annual Reviews of Computational Physics IV* (1996): 107-146.

[2] ug1399 > HLS Programmers Guide > Interface of the HLS design > Defining Interfaces > AXI Adapter Interfaces Protocol > AXI4 Master Interface > M_AXI Bundles

This hardware design has employed several hardware acceleration techniques to reduce the **latency and the initiation interval (II)**.

(1) Distributed ROM

ROM is relatively easy to replicate the storage because it only reads. The Vitis HLS tool will perform this optimization for an array that is initialized at the beginning and never modified after. Whenever the array is needed, a ROM will be synthesized, reducing the path length.

(2) Removing self-dependency

Many algorithms have to accumulate the data, however, on the hardware level, the next data can get into accumulation only after the last accumulation is finished. It is possible to **avoid such self-dependency** and increase the throughput by **alternating the logic of the C++ code** (e.g., loop interchange).

(3) Adding read ports for BRAMs

It is possible to split one large BRAM into several smaller BRAMs, increasing the number of read ports to get more data in a single clock cycle.

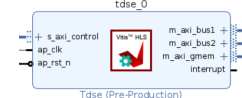
```

ps1_tmp_1_re = ps1_re[1][1];
ps1_tmp_1_im = ps1_im[1][1];
ps1_tmp_2_re = ps1_re[1+2][1];
ps1_tmp_2_im = ps1_im[1+2][1];

ps1_tmp[1][1] = ps1_tmp_1_re * c48 - ps1_tmp_2_im * s48;
ps1_im[1][1] = ps1_tmp_1_im * c48 + ps1_tmp_2_re * s48;
ps1_re[1+2][1] = ps1_tmp_1_re * s48 + ps1_tmp_2_re * c48;
ps1_im[1+2][1] = ps1_tmp_1_im * s48 + ps1_tmp_2_im * c48;

ps1_tmp_1_re = ps1_re[1][1];
ps1_tmp_1_im = ps1_im[1][1];
ps1_tmp_2_re = ps1_re[1+2][1];
ps1_tmp_2_im = ps1_im[1+2][1];

ps1_tmp[1][1] = ps1_tmp_1_re * c48 - ps1_tmp_2_im * s48;
ps1_im[1][1] = ps1_tmp_1_im * c48 + ps1_tmp_2_re * s48;
ps1_re[1+2][1] = ps1_tmp_1_re * s48 + ps1_tmp_2_re * c48;
ps1_im[1+2][1] = ps1_tmp_1_im * s48 + ps1_tmp_2_im * c48;
    
```



Removing self-dependency (left) and adding read ports (right)

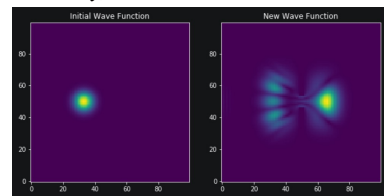
Also, it is useful to **bundle** different array arguments (with AXI4 Master interface) to different AXI4 ports if they are both accessed in one loop. This is because, by default, array arguments are mapped to a single interface bundle that only allows one read and one write in a clock cycle [2].

CREATIVE DESIGN RESULT

The hardware implementation (1.84s) of a 2-D quantum wave packet passing through an aperture is **200+ times faster** than the software (524.4s) implementation. The hardware acceleration techniques have helped reduce the initiation interval to one clock cycle, meaning a **fully pipelined** dataflow. As shown below, it only takes about 100M clock cycles to perform 2000 evolutions.

Modules && Loops	Latency/Cycles	Latency/ns	Iteration/Latency	Interval	Trip Count	Pipelined	BRAM	DSF	FF	LUT
* @ tds0	99154012	9.90068	-	99154013	-	no	134	S4	13184	16865
o sin_or_cos_float_s	18	193.000	-	1	-	yes	0	6	2497	3002
o copy_o_copy_1	10002	1.00065	4	1	10000	yes	0	6	2497	3002
* @ top	991446000	9.91068	49572	-	2000	no	-	-	-	-
o R11_o_R11_1	5014	5.01464	16	1	5000	yes	-	-	-	-
o R12_o_R12_1	4814	4.81464	16	1	4800	yes	-	-	-	-
o R21_o_R21_1	5013	5.01364	15	1	5000	yes	-	-	-	-
o R22_o_R22_1	4913	4.91364	15	1	4900	yes	-	-	-	-
o R31_o_R31_1	5012	5.01264	14	1	5000	yes	-	-	-	-
o R32_o_R32_1	4812	4.81264	14	1	4800	yes	-	-	-	-
o R41_o_R41_1	5011	5.01164	13	1	5000	yes	-	-	-	-
o R42_o_R42_1	4911	4.91164	13	1	4900	yes	-	-	-	-
o R5_o_R5_1	10053	1.01065	55	1	10000	yes	-	-	-	-

Latency, interval and resource utilization



A 2-D quantum wave packet passing through an aperture