

Towards Closing the Performance Gap for Cryptographic Kernels Between CPUs and Specialized Hardware



Naifeng Zhang



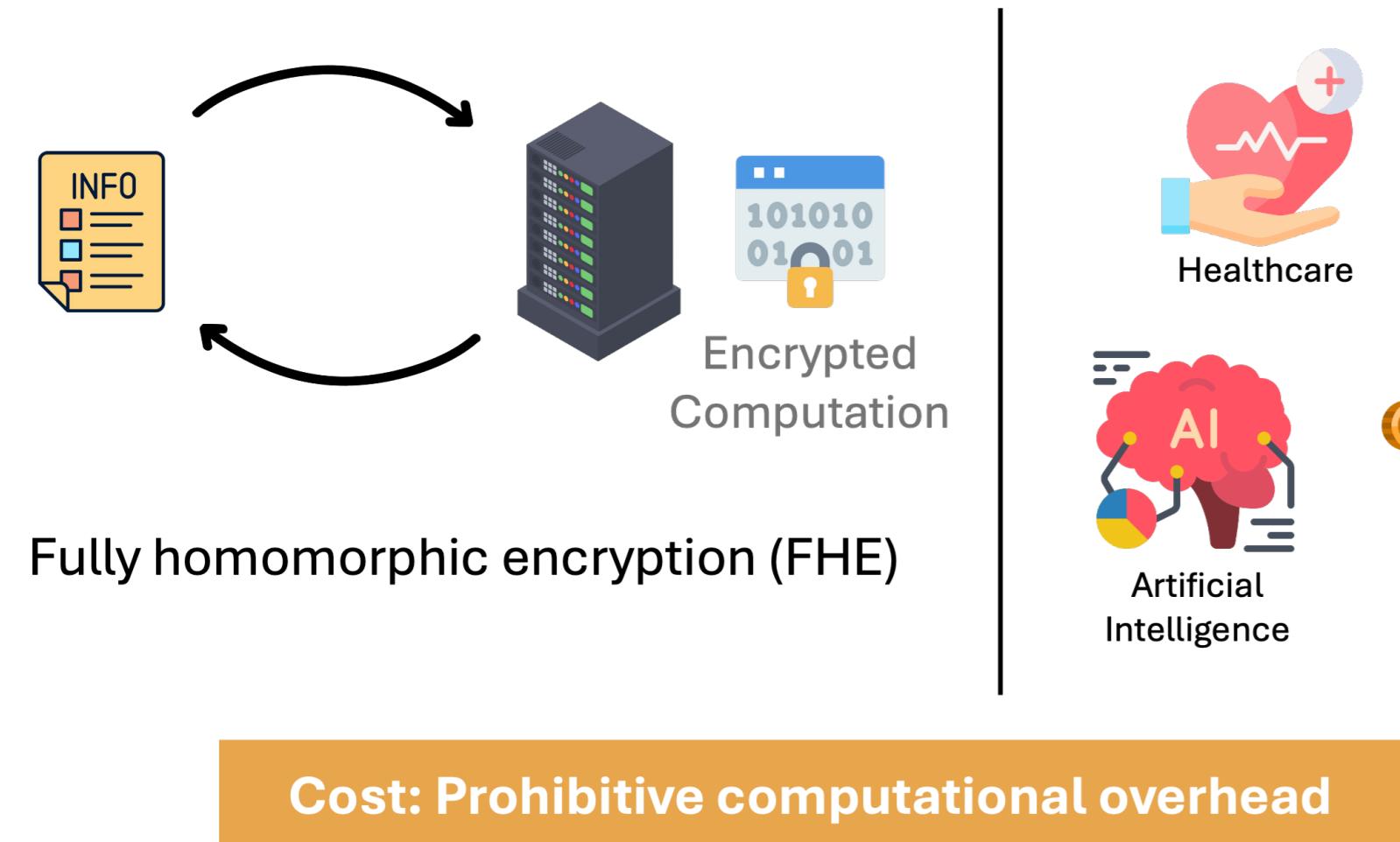
Sophia Fu



Franz Franchetti

I. MOTIVATION & BACKGROUND

Data Security Comes at a High Cost



Operations with Large Integer Arithmetic

2,000-bit integer

$$29945058656012390289752201899548961002216938026898764670213921626494125002097316037309083423551230166032468102125935343669324673615827138005859173570364330475005401605479465913340401475626944650434926155331045861080851203643664232430707675018087081999186673710699355838042146421665982617130815293083164172838554002986172170145077240914636310277685115293248536940258787378968137743709119534962195558999153247435997054900269036442663447672179509517658803138645568406212340024749768058860138217379$$

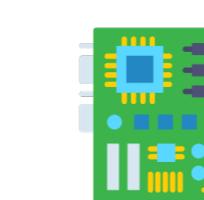
32-bit **305,419,896**

64-bit **17,920,781,457,621,970,241**

128-bit **15,107,846,090,143,992,465,023,504,163,010,990,279**

Residue Number System (RNS)

State-of-the-Art Solutions



Specialized hardware support on application-specific integrated circuits (ASICs)

- Multi-word modular arithmetic (MoMA)



Code generation-based approach on GPUs

- Multi-word modular arithmetic (MoMA)



Arbitrary precision libraries on CPUs

- GNU multiple precision (GMP) library
- OpenFHE MATHBACKEND

1,157x slowdown!

Q: HOW TO NARROW THE PERFORMANCE GAP?

II. APPROACH: SCALAR & SIMD

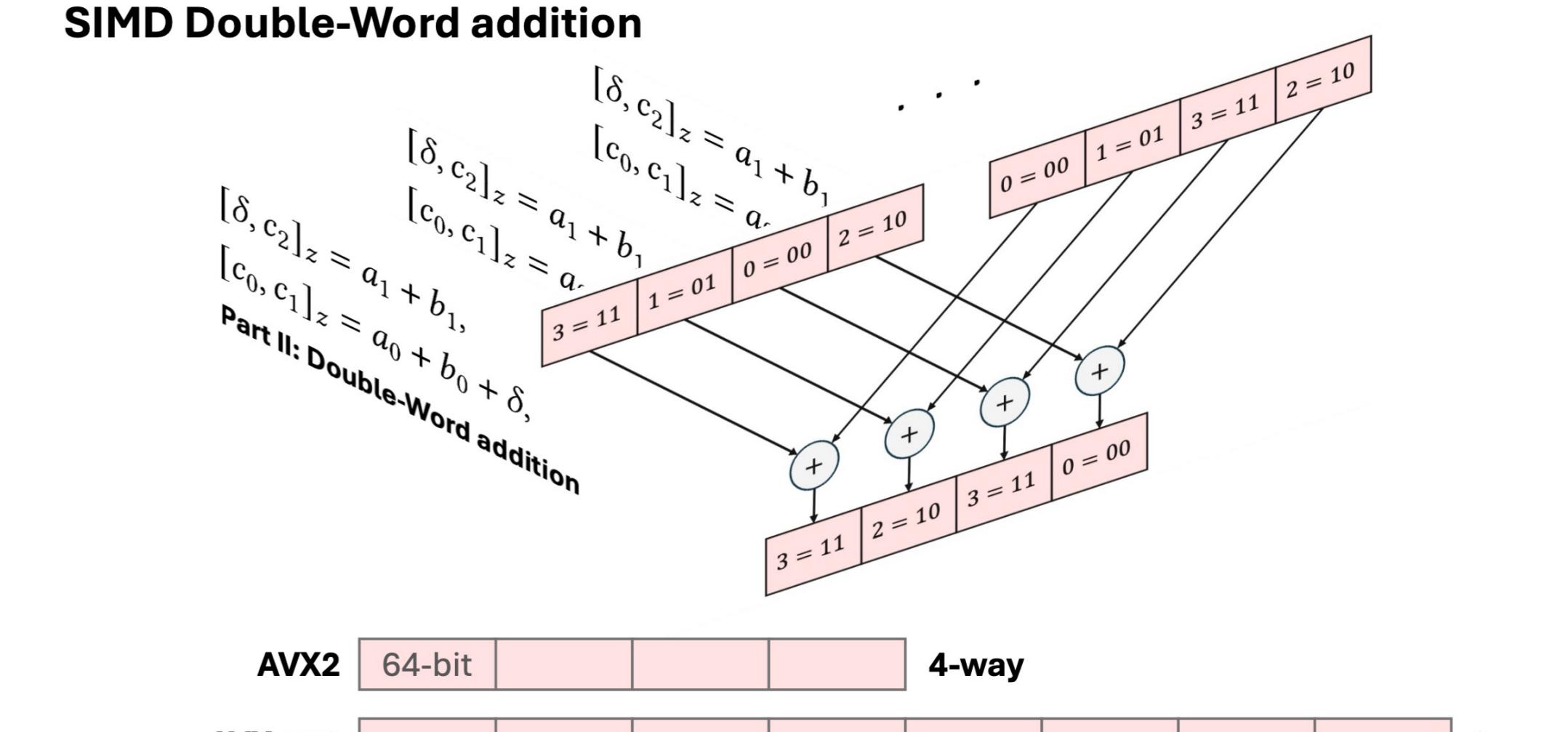
Part I: Modular Arithmetic

Math (over \mathbb{Z}_q)	Algorithm
$c = a + b \mod q$	$c = \begin{cases} a + b - q, & \text{if } (a + b) \geq q, \\ a + b, & \text{otherwise.} \end{cases}$
$c = a - b \mod q$	$c = \begin{cases} a - b + q, & \text{if } a < b, \\ a - b, & \text{otherwise.} \end{cases}$
$c = ab \mod q$	$c = ab - \lfloor ab\mu/2^k \rfloor q, \quad \mu = \lfloor 2^k/q \rfloor$ Barrett reduction

Part II: Double-Word Arithmetic

Double-word representation:	$[x_0, x_1]_z = x_0z + x_1 = x$	$z = 2^{64}$
<i>Examples</i>		
$[8, 9]_{10} = 8 \cdot 10 + 9 = 89$		
$[115292150460846975, 18446744073709550897]_{2^{64}}$ $= 21267647932558653966460912964485512497$		
Part I: Modular addition algorithm	Part II: Double-Word addition	
$c = \begin{cases} a + b - q, & \text{if } (a + b) > q, \\ a + b, & \text{otherwise.} \end{cases}$	$[\delta, c_2]_z = a_1 + b_1,$ $[c_0, c_1]_z = a_0 + b_0 + \delta,$ $\text{where } c = [c_0, c_1, c_2]_z \text{ and } \delta \in \{0, 1\}.$	
$a = [a_0, a_1]_z = a_0z + a_1$		
$b = [b_0, b_1]_z = b_0z + b_1$		

Part III: Single Instruction, Multiple Data



III. APPROACH: MULTI-WORD EXTENSION (MQX)

Instruction I: SIMD Addition with Carry

<code>_m512i_mm512_adc_epi64(_m512i a, _m512i b, __mmask8 ci, __mmask8* co)</code>									
Per-lane 64-bit addition with carry-in and outputting both the addition result and carry-out									
<i>Similar instructions</i>									
<table border="1"> <thead> <tr> <th>ISA</th><th>Instruction</th></tr> </thead> <tbody> <tr> <td>x86</td><td>ADC</td></tr> <tr> <td>LRBni</td><td>vadcp</td></tr> <tr> <td>KNC</td><td><code>_m512i_mm512_adc_epi32 (_m512i v2, __mmask16 k2, _m512i v3, __mmask16* k2_res)</code></td></tr> </tbody> </table>		ISA	Instruction	x86	ADC	LRBni	vadcp	KNC	<code>_m512i_mm512_adc_epi32 (_m512i v2, __mmask16 k2, _m512i v3, __mmask16* k2_res)</code>
ISA	Instruction								
x86	ADC								
LRBni	vadcp								
KNC	<code>_m512i_mm512_adc_epi32 (_m512i v2, __mmask16 k2, _m512i v3, __mmask16* k2_res)</code>								

Instruction II: SIMD Subtraction with Borrow

<code>_m512i_mm512_sbb_epi64(_m512i a, _m512i b, __mmask8 bi, __mmask8* bo)</code>									
Per-lane 64-bit subtraction with borrow-in and outputting the subtraction result with borrow-out									
<i>Similar instructions</i>									
<table border="1"> <thead> <tr> <th>ISA</th><th>Instruction</th></tr> </thead> <tbody> <tr> <td>x86</td><td>SBB</td></tr> <tr> <td>LRBni</td><td>vsbbpi</td></tr> <tr> <td>KNC</td><td><code>_m512i_mm512_sbb_epi32 (_m512i v2, __mmask16 k, _m512i v3, __mmask16* borrow)</code></td></tr> </tbody> </table>		ISA	Instruction	x86	SBB	LRBni	vsbbpi	KNC	<code>_m512i_mm512_sbb_epi32 (_m512i v2, __mmask16 k, _m512i v3, __mmask16* borrow)</code>
ISA	Instruction								
x86	SBB								
LRBni	vsbbpi								
KNC	<code>_m512i_mm512_sbb_epi32 (_m512i v2, __mmask16 k, _m512i v3, __mmask16* borrow)</code>								

Instruction III: SIMD Widening Multiplication

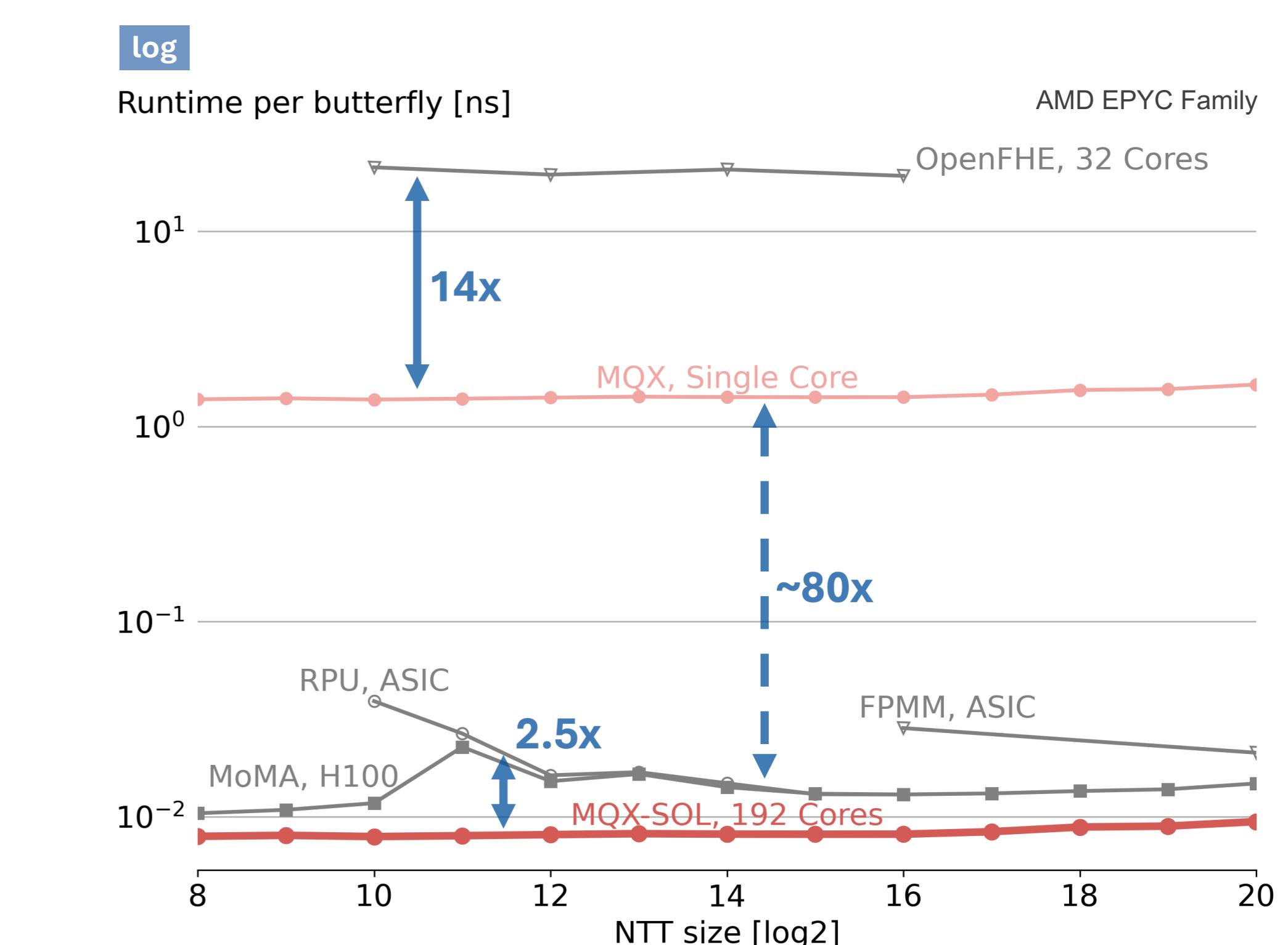
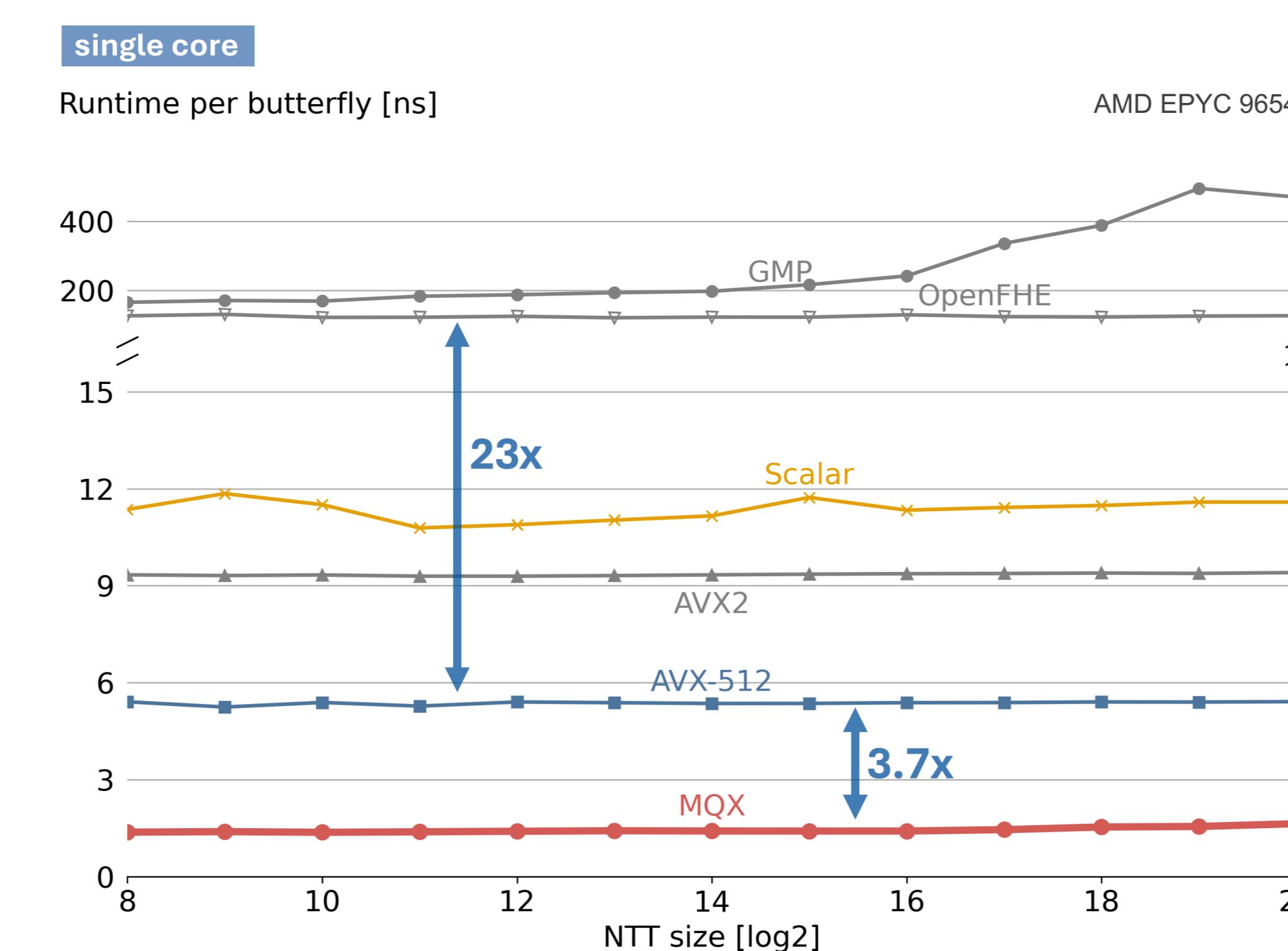
<code>void _mm512_mul_epi64(_m512i* ch, _m512i* cl, _m512i a, _m512i b)</code>											
Per each SIMD lane, multiplying two 64-bit words and storing the high part of the result in one 64-bit word and the low part in another											
<i>Similar instructions</i>											
<table border="1"> <thead> <tr> <th>ISA</th><th>Instruction</th></tr> </thead> <tbody> <tr> <td>x86</td><td>MUL</td></tr> <tr> <td>LRBni</td><td>vmulhpi</td></tr> <tr> <td>KNC</td><td><code>_m512i_mm512_mulhi_epi32 (_m512i a, _m512i b)</code></td></tr> <tr> <td>AVX-512</td><td><code>_m512i_mm512_mul_ep132 (_m512i a, _m512i b)</code></td></tr> </tbody> </table>		ISA	Instruction	x86	MUL	LRBni	vmulhpi	KNC	<code>_m512i_mm512_mulhi_epi32 (_m512i a, _m512i b)</code>	AVX-512	<code>_m512i_mm512_mul_ep132 (_m512i a, _m512i b)</code>
ISA	Instruction										
x86	MUL										
LRBni	vmulhpi										
KNC	<code>_m512i_mm512_mulhi_epi32 (_m512i a, _m512i b)</code>										
AVX-512	<code>_m512i_mm512_mul_ep132 (_m512i a, _m512i b)</code>										

IV. RESULTS

Kernel I: BLAS Operations

Kernel II: Number Theoretic Transform

Speed-of-Light Analysis



MQX ENABLES CPUs TO APPROACH ASIC-LEVEL PERFORMANCE



Code available at github.com/naifeng/benchntt
 Reach us at naifengz@cmu.edu

