

Flyte-MM: A Software Based Sub-Floating Point GEMM

Richard Veras (Louisiana State University)

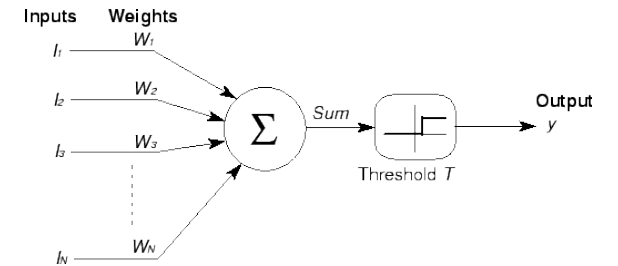
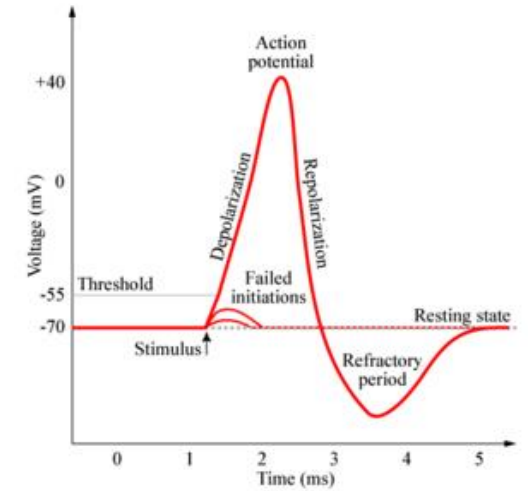
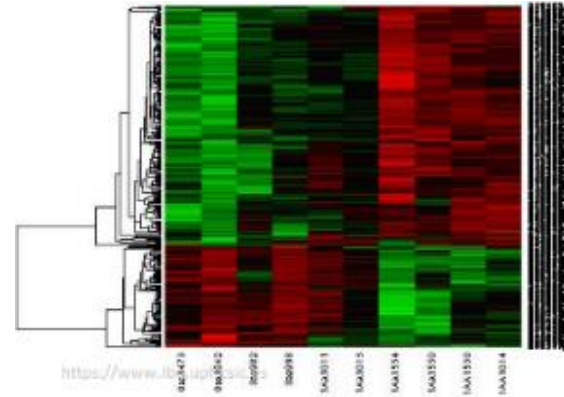
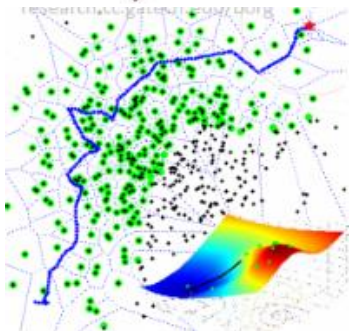
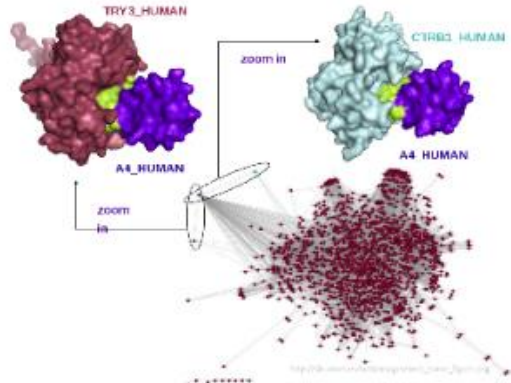
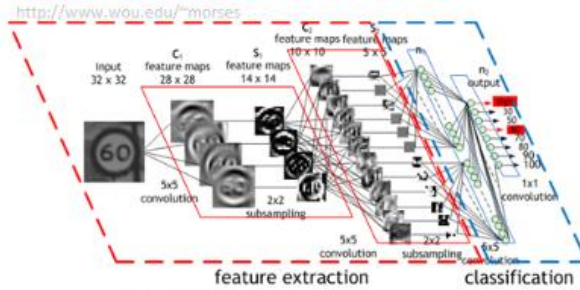
David Gregg (Trinity College, Dublin)

Want GEMM with High Dynamic Range but Low Precision

Targeting: Machine Learning and Analytics.

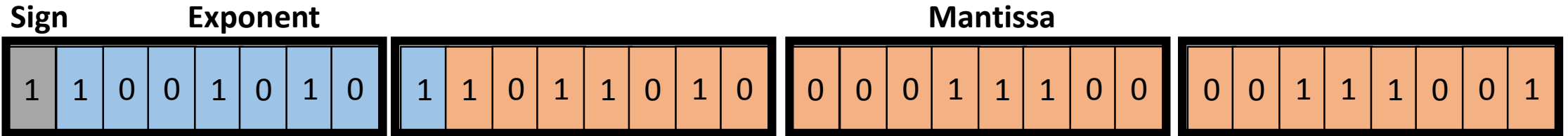
Have: Large Datasets and Constrained Systems.

Need: High Dynamic Range, Not High Precision.

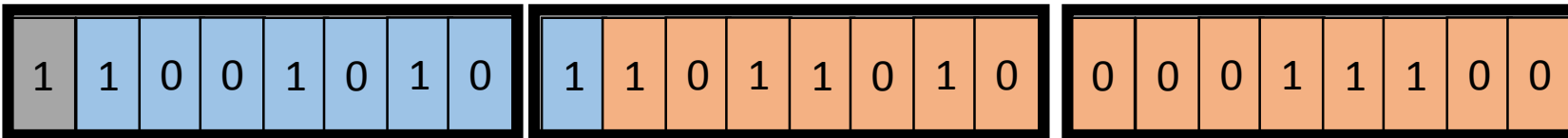


Flytes Preserve Dynamic Range at Reduced Space

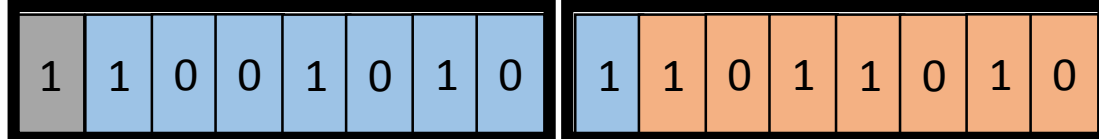
Float 32 (IEEE 754):



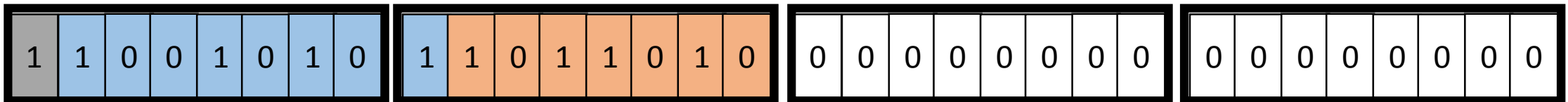
Flyte 24:



Flyte 16:

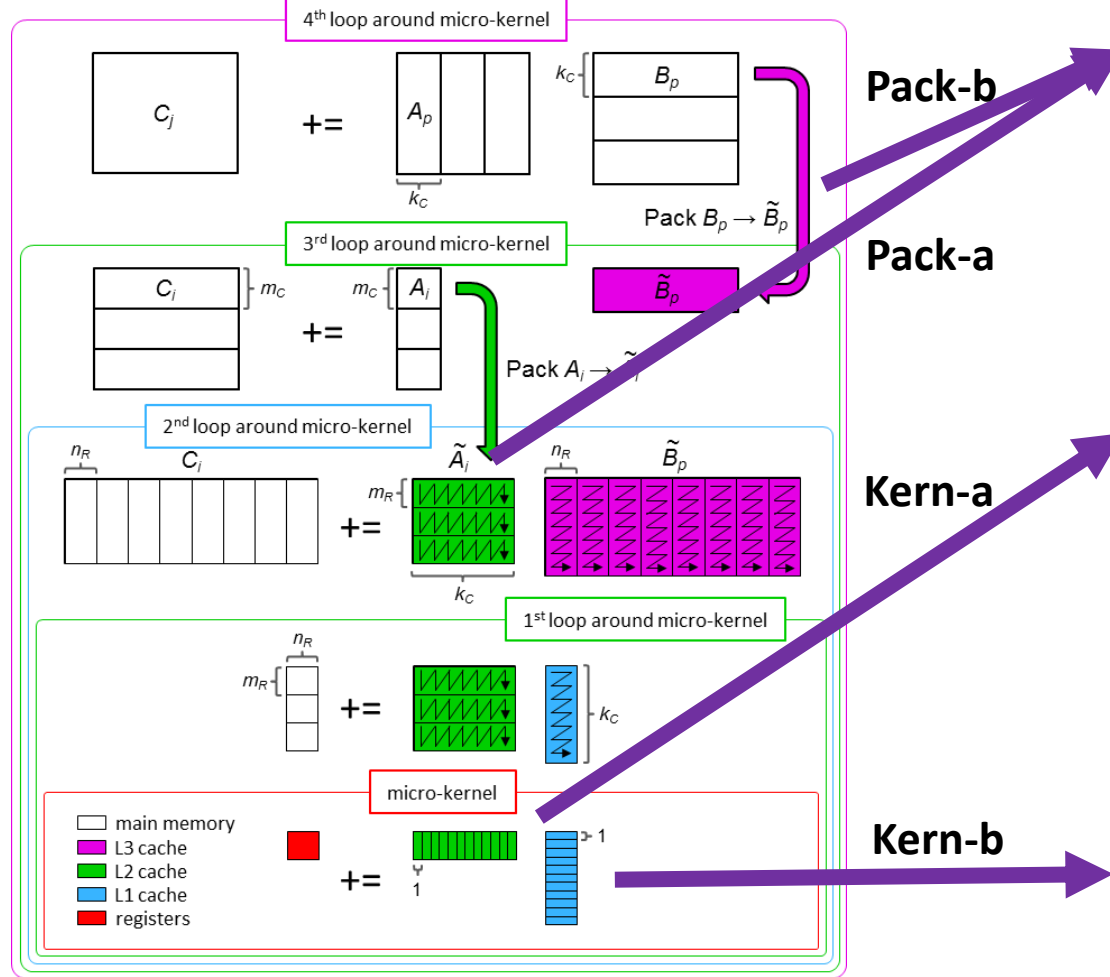


Flyte 16 to Float 32:



Efficiently Converting Floats to Flyte in GEMM

BLIS Algorithm

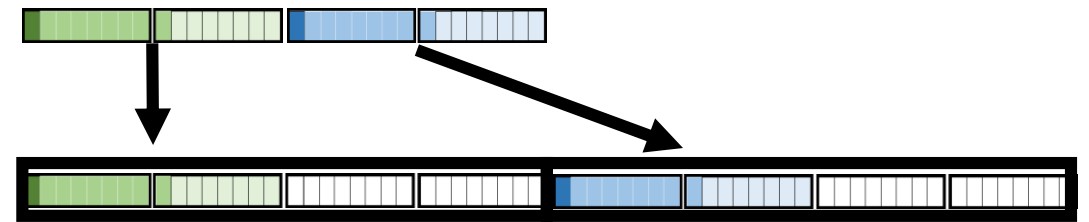


Convert in A/B in Pack: Fast, but cache block of float32.

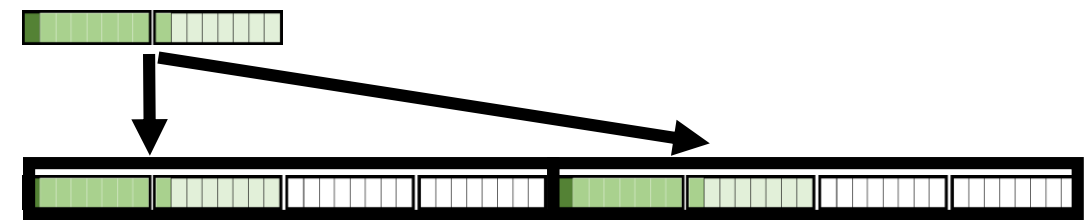
```
for( i = 0 .. mc )
  for( p = 0 .. kc )
    buffer_A[ buff_addr(i,p) ]
      = flyte_to_float(
        A[ a_addr(i,p) ] );
```

Convert in A in Kernel: Slow, but cache block of Flytes.

SIMD Load and Convert:

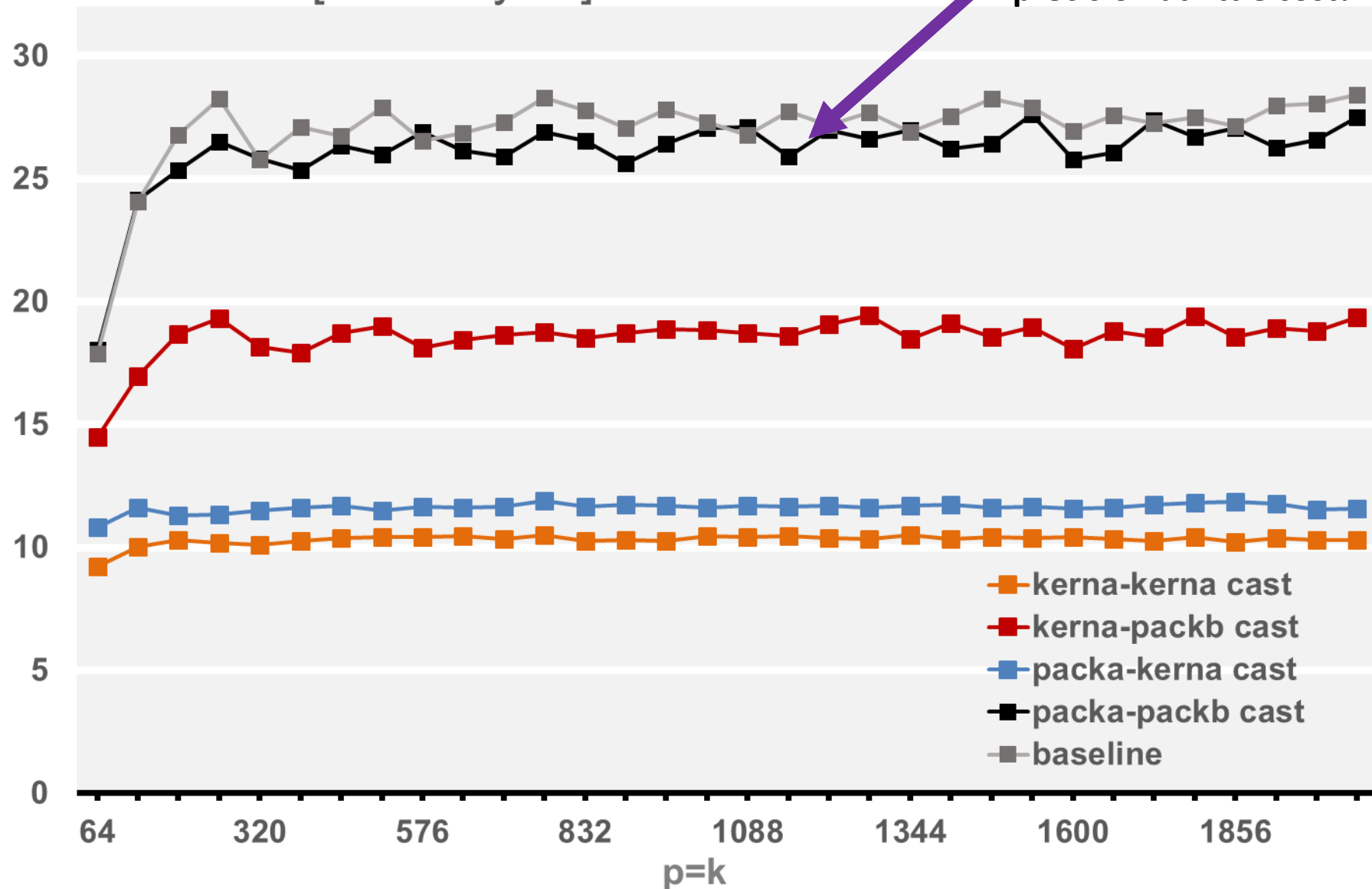


SIMD Broadcast and Convert:



A Family of Implementations

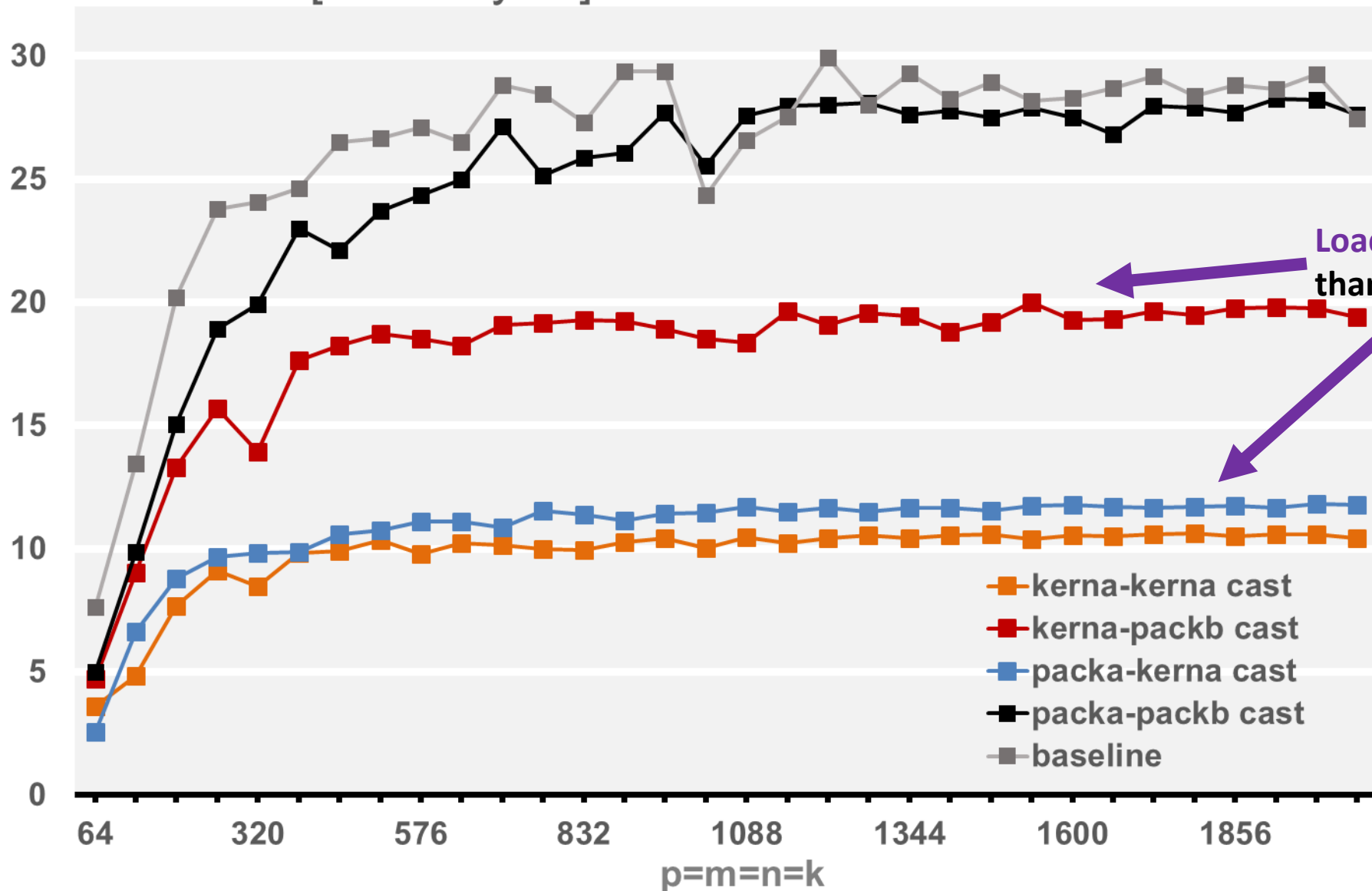
Flyte GEMM Conversion Comparison Performance [FLOP/cycle]



A Family of Implementations

Flyte GEMM Conversion Comparison

Performance [FLOP/cycle]



Load Convert is cheaper than Broadcast Convert.

Summary

- Flytes are **software based sub-precision** Floating Point.
 - Keep the sign and exponent and **truncate the mantissa**.
 - **Flyte based GEMM** desirable for **machine learning and analytics**.
- **Family of implementations** for computing on Flytes.
 - **BLIS provides** various **opportunities** to cast datatypes.
 - **Code Generation** makes **kernel** implementation **trivial**.
- Opportunity for **identifying optimal parameters**.
 - Kernel **throughput varies** with register **block dimensions**.
 - **Cache blocking** parameters **affected by datatype size**.