

In conclusion ...

- ▶ Small/skinny matrix-matrix multiplication: to pack or not to pack?
- ▶ Small MMM needs tailored algorithms
 - ▶ packing cost is not negligible
- ▶ Analyze all options
 - ▶ pack all, none, or only some matrices
 - ▶ different performance/overhead tradeoffs
 - ▶ different behaviour for skinny matrices
- ▶ In practical implementation, need for
 - ▶ good switching strategies
 - ▶ effective kernels code management
- ▶ Included in open-source BLASFEO library

