# A  Appendix

## A.1  Code

The code for benchmarks and experiments is available at `https://github.com/awav/gambit`, and the fork of TensorFlow repository with extension to XLA compiler (eXLA) is available at `https://github.com/awav/tensorflow`.

## A.2  Additional experiments

| Dataset | Distance | n | d | KeOps | eJAX | eTF | JAX | TF |
|---|---|---|---|---|---|---|---|---|
| Random | $L^2$ | 10000 | 100 | 983263 | 277364 | 284777 | 281695 | 280826 |
| Random | $L^2$ | 10000 | 10 | 2587001 | 291751 | 295029 | 287958 | 294168 |
| Random | $L^2$ | 10000 | 3 | 3662188 | 292804 | 294971 | 288098 | 294776 |
| Random | $L^2$ | 1000000 | 100 | 24367 | 2433 | 2530 | ∅ | ∅ |
| Random | $L^2$ | 1000000 | 10 | 106726 | 2505 | 2601 | ∅ | ∅ |
| Random | $L^2$ | 1000000 | 3 | 123765 | 2512 | 2605 | ∅ | ∅ |
| Random | $L^2$ | 10000000 | 100 | 2461 | 243 | 253 | ∅ | ∅ |
| Random | $L^2$ | 10000000 | 10 | 11546 | 251 | 261 | ∅ | ∅ |
| Random | $L^2$ | 10000000 | 3 | 13192 | 251 | 261 | ∅ | ∅ |
| Random | $L^1$ | 1000000 | 100 | 24307 | 517 | 521 | ∅ | ∅ |
| Random | $L^1$ | 1000000 | 10 | 108739 | 2494 | 2590 | ∅ | ∅ |
| Random | Cosine | 1000000 | 100 | 32520 | 2434 | 2515 | ∅ | ∅ |
| Random | Cosine | 1000000 | 10 | 106876 | 2507 | 2612 | ∅ | ∅ |
| MNIST | $L^2$ | 60000 | 784 | 41084 | 32290 | 33455 | 25544 | 26138 |
| MNIST | $L^1$ | 60000 | 784 | 40697 | 2356 | 2985 | 2498 | 2988 |
| Fashion | $L^2$ | 60000 | 784 | 40399 | 32382 | 33428 | 25558 | 26128 |
| Fashion | $L^1$ | 60000 | 784 | 40982 | 2357 | 2984 | 2498 | 2989 |
| Glove-50 | Cosine | 1183514 | 50 | 3464257 | 2103 | 1929 | ∅ | ∅ |
| Glove-100 | Cosine | 1183514 | 100 | 631420 | 2053 | 1871 | ∅ | ∅ |
| Glove-200 | Cosine | 1183514 | 200 | 398293 | 1967 | 1724 | ∅ | ∅ |

Table 3: Query processing rates (queries per second) for kNN. $n$ and $d$ are the number of data points and the data dimension respectively. Runs which failed due to memory overflow are denoted by ∅. Runs with eXLA are denoted eJAX and eTF respectively.

Figures 5 and 6 depict XLA HLO graphs for kernel matrix-vector multiplication before and after splitting optimisation in eXLA optimisation pipeline (section 4.3), respectively. The same configuration of the kernel is used as in section 5.1, i.e. squared exponential kernel from Matthews et al. (2017). By kernel matrix-vector multiplication expression we mean the function $g(\mathbf{x}, \mathbf{y}, \mathbf{v}) = k(\mathbf{x}, \mathbf{y})\mathbf{v}$, where $k(\mathbf{x}, \mathbf{y}) = \sigma^2 \exp\left(-1/2\|\mathbf{x} - \mathbf{y}\|^2/l^2\right)$ is the kernel with $\sigma^2$ and $l$ hyperparameters. The size of 1-dimensional input vectors $\mathbf{v}$, $\mathbf{x}$ and $\mathbf{y}$ is $1e{-}6$. In turn, the size of the corresponding kernel matrix is $1e{-}6 \times 1e{-}6$, and in double precision would require to allocate 8TB. The *tensor size threshold* was set to 1GB, and eXLA splitting optimisation pass divided the expression of the kernel matrix-vector multiplication into smaller chunks, such that the maximum tensor size in the graph is $125 \times 1e{-}6$.

Figure 5: XLA HLO graph for kernel matrix-vector multiplication **before** splitting optimisation pass is applied in the XLA optimisation pipeline.
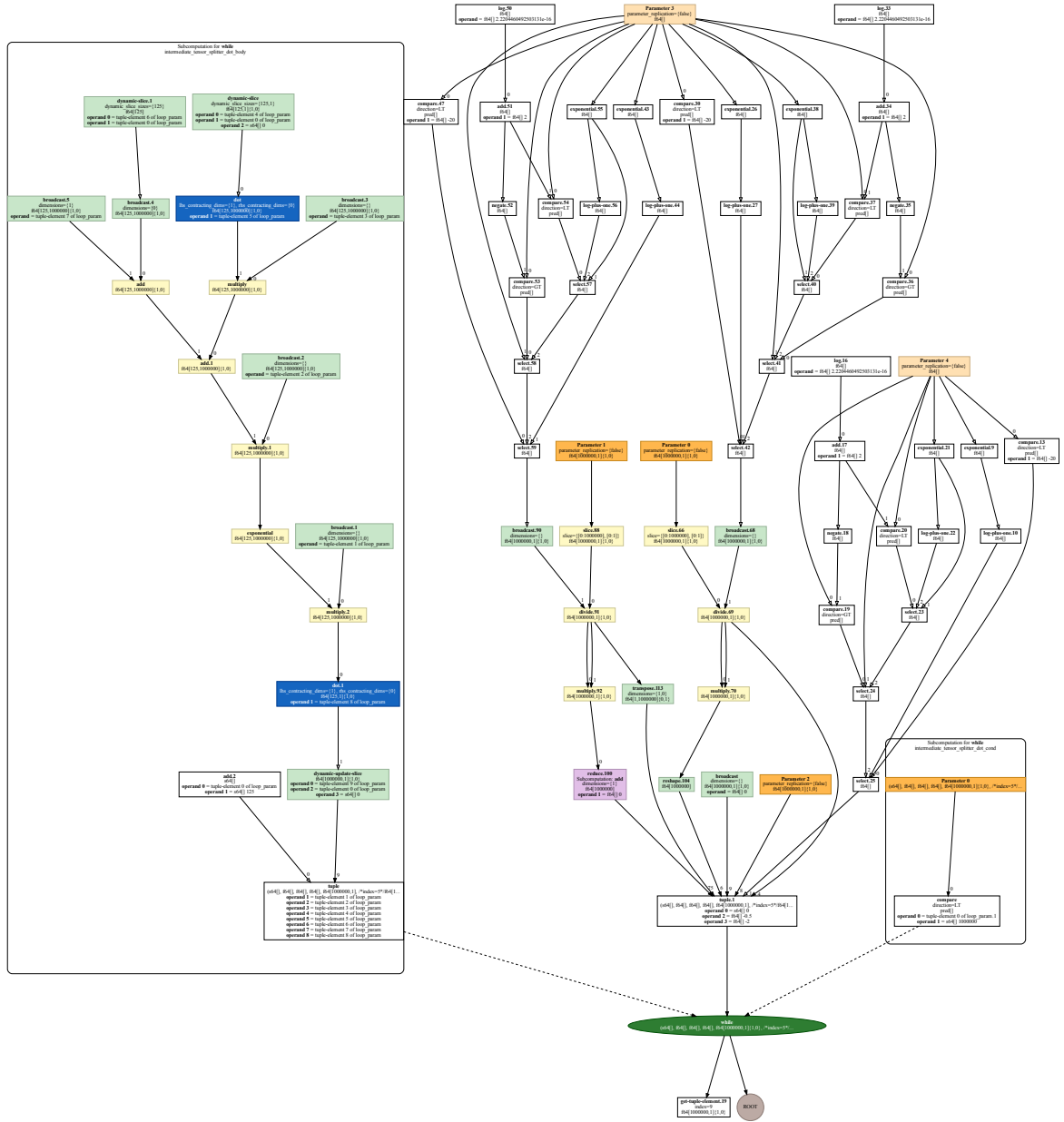
Figure 6: XLA HLO graph for kernel matrix-vector multiplication **after** splitting optimisation pass is applied in the XLA optimisation pipeline.