



LAWRENCE
LIVERMORE
NATIONAL
LABORATORY

Data Placement Optimization in GPU Memory Hierarchy using Predictive Modeling

L. Stoltzfus, M. Emani, P. Lin, C. Liao

September 11, 2018

Workshop on Memory Centric High Performance Computing
Dallas, TX, United States
November 11, 2018 through November 11, 2018

Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

Data Placement Optimization in GPU Memory Hierarchy using Predictive Modeling

Larisa Stoltzfus
University of Edinburgh
Edinburgh, UK
larisa.stoltzfus@ed.ac.uk

Pei-Hung Lin
Lawrence Livermore National Laboratory
Livermore, CA
lin32@llnl.gov

Murali Emami
Lawrence Livermore National Laboratory
Livermore, CA
emami1@llnl.gov

Chunhua Liao
Lawrence Livermore National Laboratory
Livermore, CA
liao6@llnl.gov

ABSTRACT

Modern supercomputers often use Graphic Processing Units (or GPUs) to meet the ever-growing demands for high performance computing. GPUs typically have a complex memory architecture with various types of memories and caches, such as global memory, shared memory, constant memory, and texture memory. The placement of data on these memories has a tremendous impact on the performance of the HPC applications and identifying the optimal placement location is non-trivial.

In this paper, we propose a machine learning-based approach to build a classifier to determine the best class of GPU memory that will minimize GPU kernel execution time. This approach utilizes a set of performance counters obtained from profiling runs and combines with relevant hardware features to generate the trained model. We evaluate our approach on several generations of NVIDIA GPUs, including Kepler, Maxwell, Pascal, and Volta on a set of benchmarks. The results show that the trained model achieves prediction accuracy over 90% and given a global version, the classifier can accurately determine which data placement variant would yield the best performance.

KEYWORDS

GPU, Data placement, Memory, Machine Learning, Predictive Modeling

1 INTRODUCTION

Supercomputers use various types of memory and increasingly complex designs to meet the ever growing needs for data by modern massively parallel processors (e.g., Graphic Processing Units or GPUs). On an NVIDIA Kepler GPU, for instance, there are more than eight types of memory (global, texture, shared, constant, various cache, etc.). Figure 1 shows a typical GPU memory hierarchy. With unified CPU-GPU memory space support on the latest GPUs (e.g., NVIDIA Volta), data placement becomes even more flexible, allowing direct accesses to data across the boundaries of multiple

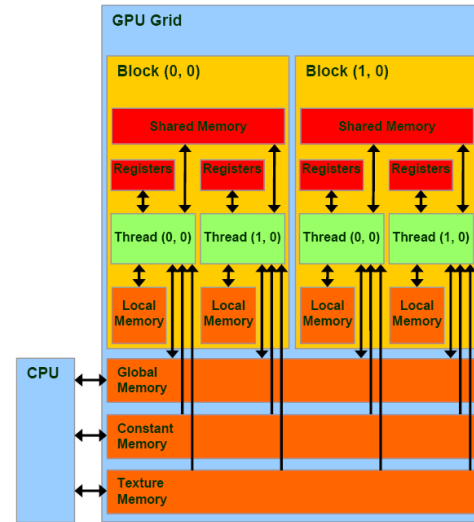


Figure 1: GPU memory hierarchy

GPUs and CPU. This problem will critically affect the effective adoptions of new generations of supercomputers, such as Sierra hosted at LLNL [1] and Summit at ORNL [2] featuring NVIDIA Volta GPUs using 3D stacked memory, unified CPU-GPU memory space, and other memory complexities. As HPC applications get ported to run on the new supercomputers with GPUs, they will have to rely on a range of optimizations including data placement optimization to reach desired performance.

Studies [4, 8] have shown that placing data onto the proper part of a memory system called *data placement optimization*, has significant impact on program performance; it is able to frequently speed up carefully written GPU kernels by over 50% (up to 13.5× on an MPEG video encoding kernel [7]). Previous studies have explored various approaches of optimizing data placement on early generations of GPUs. For example, PORPLE [3, 5] is a portable GPU data placement approach combining a memory specification, a compiler framework and a run-time library. A lightweight performance model based on reuse distance is used by PORPLE to guide run-time selection of optimal data placement policies. Its effectiveness has been evaluated on Tesla, Fermi, and Kepler GPUs. Huang and Li [6]

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under contract DE-AC52-07NA27344 and was supported by the LLNL-LDRD Program under Project No. 18-ERD-006. (LLNL-CONF-XXXXX).

have proposed a new performance modeling approach for optimal data placement on GPUs. They have analyzed correlations among different data placements and used a sample data placement to predict performance for other data placements. The work has been evaluated using a Kepler GPU. Jang et al. [7] have presented several rules based on data access patterns to guide the memory selection for a Tesla GPU. Mei et al. [11] have proposed a new fine-grained microbenchmarking approach to expose unknown cache parameters of three generations of GPUs, including Fermi, Kelper, and Maxwell. More recently, Jia et al. [9] have used microbenchmarking to analyze the Volta GPU architecture.

The novelty in this work is in predicting optimal data placement on newer GPUs. Prior state-of-art approaches used offline profiling analysis but are unable to capture the required features at runtime. By observing the runtime features of just the global/default version of a code, the proposed approach is able to determine which memory placement will yield best performance. In this work, we use machine learning based model that is trained once offline. During inference, runtime features are captured and passed as an input feature vector to the model that determines the optimal data placement location.

This paper makes the following novel contributions:

- determining the optimal data placement location on-the-fly during run-time,
- we propose a lightweight solution provide a simple resulting solution that is applicable to diverse applications
- the approach and data can be reused for other optimizations such as determining optimal data layouts.

2 MOTIVATION

In this section, we demonstrate how the data placement in memory could significantly impact the program execution time. Here we first provide a brief description of GPU memory hierarchy and then show the impact of different placement of data onto various types of memories on the kernel execution time.

2.1 GPU Memory Structure

GPUs have a highly complex memory hierarchy in order to exploit their massive parallel computing potentials. Traditionally they are characterized by a hierarchy of memory levels, namely: global, constant, texture, shared and registers. A high-level overview of the major types of memories as exposed to programmers via the CUDA API is listed as follows:

- **Global Memory:** This largest off-chip memory serves as default, main memory on a GPU. Global memory accesses are characterized with limited memory bandwidth and long latencies compared with on-chip memory or cache.
- **Shared Memory:** This on-chip memory is characterized by low-latency and high-bandwidth. It is software-managed and is accessible by active threads that belong to a streaming multiprocessor (SM) unit on a GPU.
- **Constant Memory:** This is a predefined within the global memory space and is set as read-only. The memory space is globally visible to all threads.
- **Texture Memory:** Similar to read-only constant memory, texture memory is off-chip memory space that is optimized

Larisa Stoltzfus, Murali Emani, Pei-Hung Lin, and Chunhua Liao

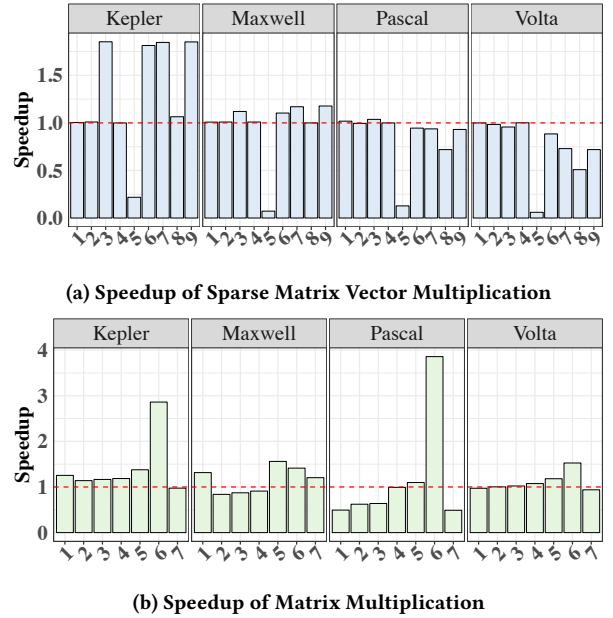


Figure 2: Impact of data placement on kernel performance across different GPUs. Speedups over default vary with different memory variants shown on x-axis.

for 2D spatial locality. It is suited for threads that access memory addresses that are closer to each other in 2D.

Although all NVIDIA GPUs have similar high-level design, different generations of GPUs introduce new memory properties or implement the same memories using different physical organizations. For example, Fermi GPUs introduced true cache hierarchy for global memory while previous GPUs did not have such caches. On Kepler GPUs, L1 cache and shared memory are combined together while texture cache is backed by its own memory chip. Pascal and Maxwell GPUs have dedicated memory for their shared memory. For the latest Volta GPUs, L1 cache, shared memory and texture cache are all merged into a single unified memory. All these hardware changes have direct impact of memory optimization.

2.2 Preliminary Experiments

Preliminary experiments show that the performance achievable from different data placements can be specific to applications and platforms. From Figure 2, it is evident that the version of the hardware and the application, both heavily influence the kernel performance. The graphs show benchmarks utilizing different GPU memories, numbered across the x-axis and their speedup compared to the global version of the benchmark. For example, Table 1 shows what different numbered versions correspond to for the SPMV benchmark (for brevity, only one table is shown). Figure 2a shows that the performance of SPMV application, which is a sparse matrix vector application, generally worsens with the latest version of NVIDIA GPUs. However for MM, the matrix-matrix multiplication benchmark, performance improvement can still be seen on the latest platforms (Volta and Pascal) as shown in Figure 2b. Both

benchmarks show that for many applications it is clearly important to get the data placement right. Applications optimized for a particular platform will not necessarily retain that performance on future platforms and indeed future platforms may not require optimizing at all.

Version	Description
1	rows array in shared
2	rows array in constant
3	vector array in texture1D, rows in shared
4	matrix values in texture1D
5	vector array in constant, rows in texture1D
6	vector array in texture
7	matrix values and columns in texture1D, rows in constant
8	matrix values, columns and vector in texture1D
9	matrix values, columns, rows and vector in texture1D

Table 1: Details of SPMV Memory Configurations

3 APPROACH

The machine learning model for directly predicting the appropriate placement plan is based on the access patterns and the memory system specifications. The workflow involves two phases, namely offline training and online inference, as shown in Figure 3. The offline training phase involves the construction of a classifier using supervised learning techniques. It takes the feature vector of the kernel and data as inputs, and predicts the best data placement in memory. As with any machine learning approach, we investigate different classifiers (e.g., Random Forests, Support Vector Machines, AttributedClassifiers) and pick one that yields the highest prediction accuracy and confidence. To train the machine learning models, we obtain the training data by evaluating representative kernels with possible placements and labeling the samples with observed best execution times.

3.1 Offline Training

3.1.1 Design of Training Experiments. The training experiments are set up to run scripts over a selection of benchmarks with varying types of data placements and data and thread block sizes across different GPU platforms listed in Section 4. First, the best performing versions are found by measuring the program execution times of all possible memory variants, for an array of data sizes and thread block sizes of each benchmark. Then the global memory versions are profiled for their feature values, selecting only those metrics and events available on all platforms. These features are paired with class labels of the best versions available for a given benchmark of a data size and thread block size. Feature values are first normalized and then re-sampled to ensure a larger spread for use. The benchmarks along with the hardware platforms can be found in Section 4.

3.1.2 Feature Engineering: As with any machine learning-based models, the classifier to make accurate predictions for the data placement, relies heavily on the input features. Here, we present how the features are extracted and important ones are selected.

(a) *Feature Extraction:* We use the NVIDIA profiling tool nvprof to extract the hardware counters and performance metrics that

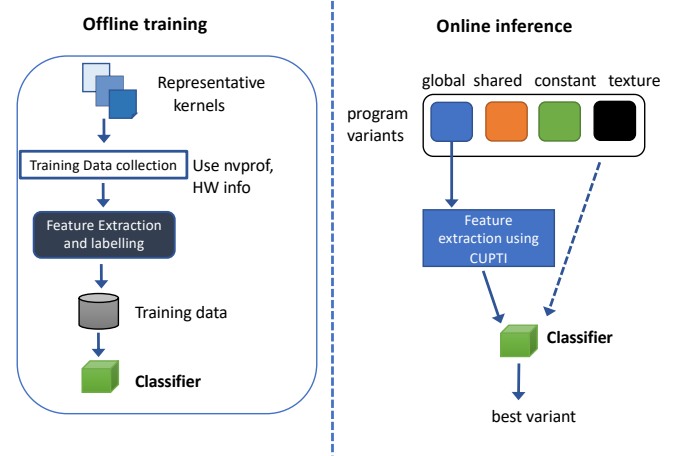


Figure 3: Workflow of the proposed machine learning based data placement engine

Feature	Description
<i>active_warps</i>	number of active warps
<i>gld_transactions</i>	average number of global memory load transactions performed for each global memory load
<i>stall_exec_dependency</i>	percentage of stalls occurring because an input required by the instruction is not yet available
<i>achieved_occupancy</i>	ratio of the average active warps per active cycle to the maximum number of warps
<i>ldst_executed</i>	Number of executed local, global, shared and texture memory load and store instructions

Table 2: List of selected features to train the classifier

compose the features and use the CUPTI API [12] to extract these features from the codes at runtime. Hardware parameters obtained from extensive microbenchmarking [10] are then appended to each sample of features obtained nvprof. The raw training data set is obtained by running 9 benchmark programs with various combinations of 4 memory variants, utilizing 3 variations of data sizes and 3 different thread block sizes on 4 GPU architectures. This results in a total of 1296 samples with 241 features per sample collected from nvprof and hardware.

(b) *Dimensionality Reduction:* The dimensionality of the features is reduced from 241 to 5 (Table 2); these significant ones are obtained using the correlation based feature selection (CFS) algorithm. Such reduction in the required features helps to extract only those features that are important to the model.

3.1.3 Model Training: With the reduced feature set, we evaluated different classifiers to compare respective prediction accuracies, obtained using 10-fold cross validation. Such a strategy ensures that the model does not overfit the training data and the reported

Classifier	Prediction accuracy
RandomForest	95.7%
LogitBoost	95.5%
IterativeClassifierOptimizer	95.5%
SimpleLogistic	95.4%
JRip	95.0%

Table 3: Classifiers and their prediction accuracies

prediction accuracies are based on unseen test data. Classifiers with high prediction accuracies are listed in Table 3.

3.2 Online Inference

The classifier built in the offline training phase is then used for predicting the best data placement for new applications at runtime through the CUPTI profiler API. Honing in on the fewest number of features needed, the application is profiled in real-time, the features are then fed as input to the classifier. Based on the prediction, the best version of the application is then executed. The assumption here is that different memory variants of the code exist.

4 EVALUATION

Here we list the hardware and software platform used to evaluate the proposed approach, followed by the experimental results that show how close the model predicted memory variant is with the best possible one. Table 4 show four machines for our experiments. They contain four generations of GPUs namely Kepler, Maxwell, Pascal and Volta. The programs include *Sparse Matrix Vector Multiplication* (SPMV), *Molecular Dynamics Simulation* (MD), *Computational Fluid Dynamics Simulation* (CFD), *Matrix-Matrix Multiplication* (MM), *ParticleFilter*, *ConvolutionSeparable*, *Stencil27*, and *Trans*.

In the first run, the best performing version of each benchmark is determined. The benchmarks were run using a script to collect average values over ten iterations and the median of these values is taken. Each kernel was run for five times to warm up the GPU before timings were taken. In order to collect consistent performance times, the performance benchmarks used `cudaEventRecord()` to denote kernel start and stop places and `cudaEventElapsedTime()` to calculate the time elapsed. `cudaDeviceSynchronize()` was used after these calls and data transfer times are excluded in order to isolate performance differences from memory usage.

5 ANALYSIS

The results of our machine learning model experiments show that we are able to get very good results from tree-based models, however we are not quite able to replicate this level of accuracy when applying the model to a new application. Below we show two graphs depicting these two results. First we show how well the JRIP model performs and the accuracy in its predictions for each of the benchmarks across the different platforms. We have selected this model because it is the most portable of those available and all of the top performing models hover around ~95% prediction accuracies. Next we discuss how well this model performs on a new application: a room acoustics model benchmark. Even if the model

may not accurately predict the best performing data placement, it may still predict a better data placement than the global version.

The graph in Figure 4 shows the results for the best data placement class comparing the model predicted memory variants with the best performing ones. We found that the memory class predicted by the machine learning model is accurate in at least 80% and up to 95% of cases. Specially, given a global memory variant, this model is able to correctly predict which of the memory variants would yield the best performance (lowest execution time). These results are averaged across different GPUs. It can be observed from the figure that, only in few circumstances, the model misclassifies the best memory variant. Even here, the difference in the execution time was found to no more than 22%.

In order to prove that our prediction accuracy is not just limited to the benchmarks considered, we tested our model on a room acoustics benchmark[13] which does not exist in the training set. As shown in Figure 5, the machine learning based model is able to correctly predict the memory variant that has the best performance. As in earlier case, it has particularly higher accuracy for texture and global than shared and constant memory variants. This result demonstrates that the proposed model is portable and is able to determine what memory placement variant minimizes the execution times with totally new programs not included in the training set.

6 CONCLUSION

We have presented an automated approach to data placement optimization using machine learning to tune applications on GPUs. We built a classifier to determine the memory variant that could yield highest performance. The evaluations demonstrate that the model predictions are nearly accurate as the best achievable cases.

This work has shown that there is an immense possibility for machine learning to be applied to automate data placement optimizations. However, the long term goal would be to have this integrated into an automated workflow and have runtime code generation of the appropriate memory variant based on the model decision. It would be ideal for a compiler to handle these low-level details. Additionally, the overhead of using CUPTI means that performance suffers. Essentially, the profiler must run the kernel of interest for each metric or event, meaning that extra iterations are added. This could be alleviated by running cut down data sizes. Additionally, this work could easily be applied to other areas such as different data layout optimizations.

REFERENCES

- [1] 2018. Sierra- next generation HPC system at LLNL. <https://computation.llnl.gov/computers/sierra>. (2018).
- [2] 2018. Summit at Oak Ridge Leadership Computing Facility. <https://www.olcf.ornl.gov/summit/>. (2018).
- [3] Guoyang Chen, Xipeng Shen, Bo Wu, and Dong Li. 2017. Optimizing data placement on GPU memory: A portable approach. *IEEE Trans. Comput.* 66, 3 (2017), 473–487.
- [4] Guoyang Chen, Bo Wu, Dong Li, and Xipeng Shen. 2014. PORPLE: An Extensible Optimizer for Portable Data Placement on GPU. In *The 47th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. 88–100.
- [5] Guoyang Chen, Bo Wu, Dong Li, and Xipeng Shen. 2014. PORPLE: An extensible optimizer for portable data placement on GPU. In *Proceedings of the 47th Annual IEEE/ACM International Symposium on Microarchitecture*. IEEE Computer Society, 88–100.
- [6] Yingchao Huang and Dong Li. 2017. Performance modeling for optimal data placement on GPU with heterogeneous memory systems. In *Cluster Computing (CLUSTER), 2017 IEEE International Conference on*. IEEE, 166–177.

Data Placement Optimization in GPU Memory Hierarchy using Predictive Modeling

	Kepler (K40)	Maxwell (M60)	Pascal (P100)	Volta (V100)
CPU	IBM Power8 @2.2GHz	Intel Xeon E5-2670 @2.60 GHZ	IBM Power8 @2.2GHz	Intel Xeon E5-2699 @2.20GHz
Computation capability	3.5	5.2	6.0	7.0
SMs	15	16	56	80
Cores/SM	192 SP cores/64 DP cores	128 cores	64 cores	64 SP cores/32 DP cores
Texture Units/SM	16	8	4	4
Register File Size/SM	256 KB	256 KB	256 KB	256 KB
L1 Cache/SM	Combined L1+Shared 64K	Combined 24KB	Combined 24 KB	128 KB Unified
Texture Cache	48KB			
Shared Memory/SM	Combined L1+Shared 64K			
L2 Cache	1536 KB	2048 KB	4096 KB	6144KB
Constant Memory	64 KB	64 KB	64 KB	64 KB
Global Memory	12 GB	8 GB	16 GB	16 GB

Table 4: Key specifications of selected GPUs of different generations

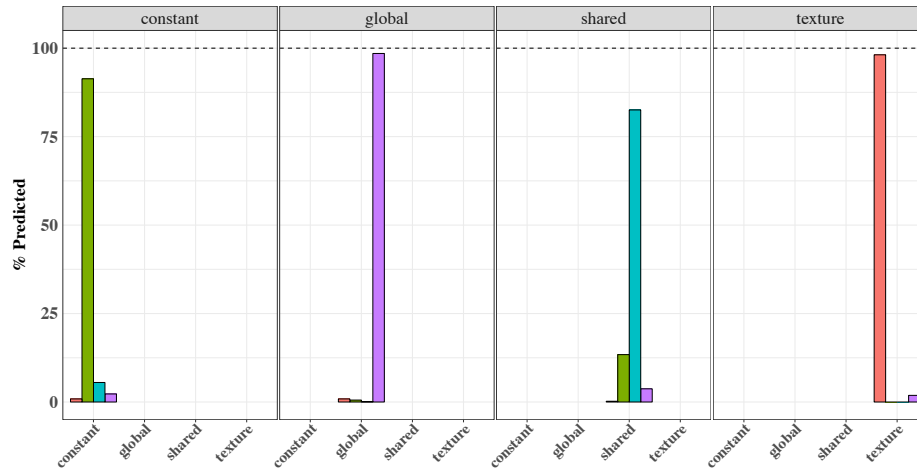


Figure 4: Graph showing the model-predicted memory classes with the best performing memory variants. Across all GPUs the model correctly predicts the best performing variant in at least 80% and up to 95% cases.

- [7] Byunghyun Jang, Dana Schaa, Perhaad Mistry, and David Kaeli. 2010. Exploiting memory access patterns to improve memory performance in data-parallel architectures. *IEEE Transactions on Parallel & Distributed Systems* 1 (2010), 105–118.
- [8] Keon Jang, Sangjin Han, Seungyeop Han, Sue Moon, and KyoungSoo Park. 2011. SSLShader: Cheap SSL Acceleration with Commodity Processors. In *The 8th USENIX Conference on Networked Systems Design and Implementation (NSDI)*. 1–14.
- [9] Zhe Jia, Marco Maggioni, Benjamin Staiger, and Daniele P Scarpazza. 2018. Dissecting the NVIDIA Volta GPU Architecture via Microbenchmarking. *arXiv preprint arXiv:1804.06826* (2018).
- [10] Zhe Jia, Marco Maggioni, Benjamin Staiger, and Daniele Paolo Scarpazza. 2018. Dissecting the NVIDIA Volta GPU Architecture via Microbenchmarking. *CoRR* abs/1804.06826 (2018). [arXiv:1804.06826](http://arxiv.org/abs/1804.06826) <http://arxiv.org/abs/1804.06826>
- [11] Xinxin Mei and Xiaowen Chu. 2017. Dissecting GPU memory hierarchy through microbenchmarking. *IEEE Transactions on Parallel and Distributed Systems* 28, 1 (2017), 72–86.
- [12] NVIDIA. 2018. CUDA Profiling Tools Interface. (2018). <https://developer.nvidia.com/cuda-profiling-tools-interface>
- [13] Craig Webb. 2014. *Parallel Computation Techniques for Virtual Acoustics and Physical Modelling Synthesis*. Ph.D. Dissertation. University of Edinburgh. www.ness-music.eu/wp-content/uploads/2014/07/CJWebb_thesis-1.pdf

[ness-music.eu/wp-content/uploads/2014/07/CJWebb_thesis-1.pdf](http://www.ness-music.eu/wp-content/uploads/2014/07/CJWebb_thesis-1.pdf)

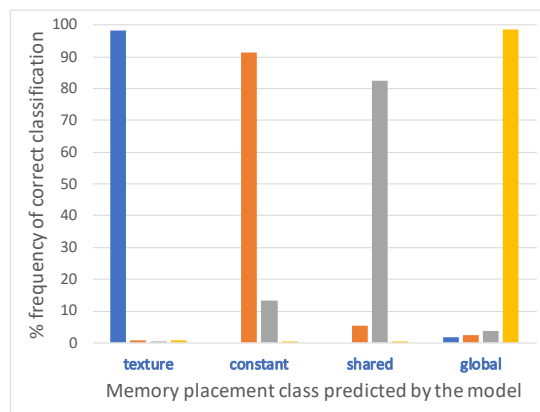


Figure 5: Predictions for Acoustic Benchmark