

NC State University
Department of Electrical and Computer Engineering
ECE 463/563: Fall 2019 (Dr. Huiyang Zhou)
Project #2: Branch Prediction

By
Rhushikesh Anand Prabhune

NCSU Honor Pledge: “I have neither given nor received unauthorized aid on this test or assignment.”

Student’s Electronic Signature: Rhushikesh Anand Prabhune
(Sign by typing your name)

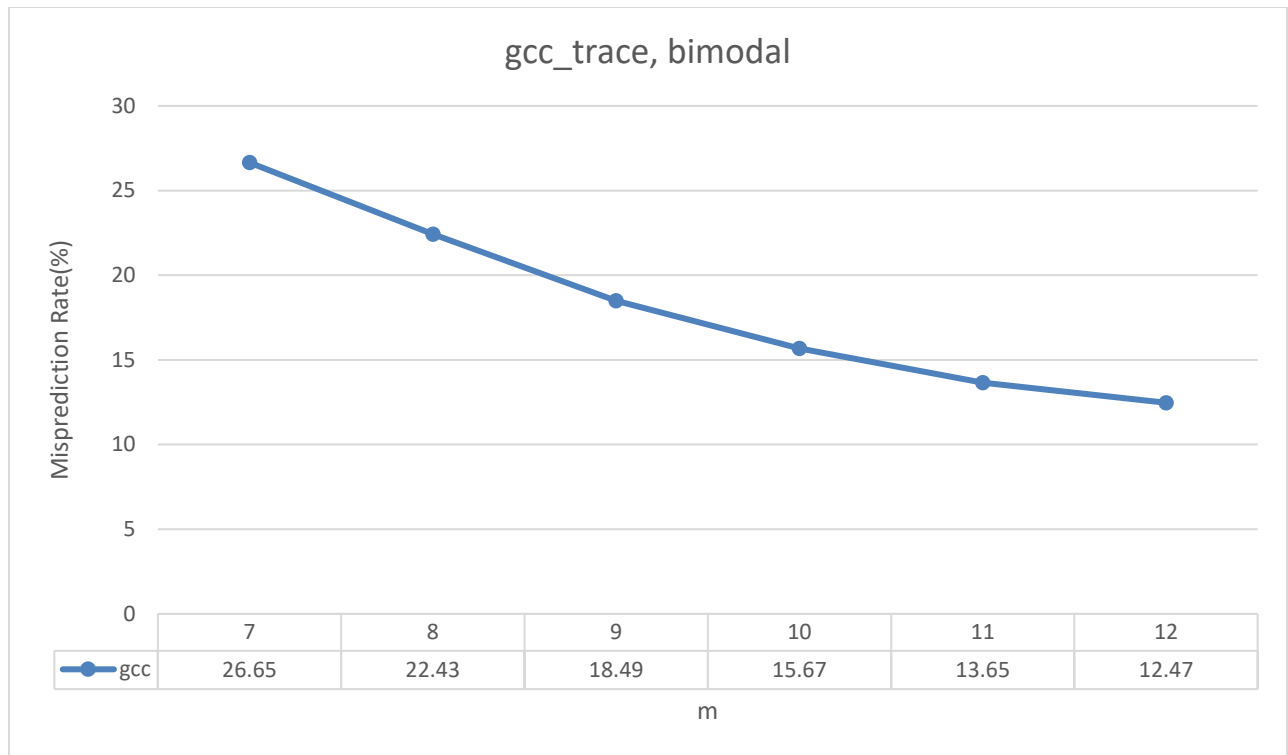
Course Number: 563
(463 or 563?)

Bimodal Predictor:

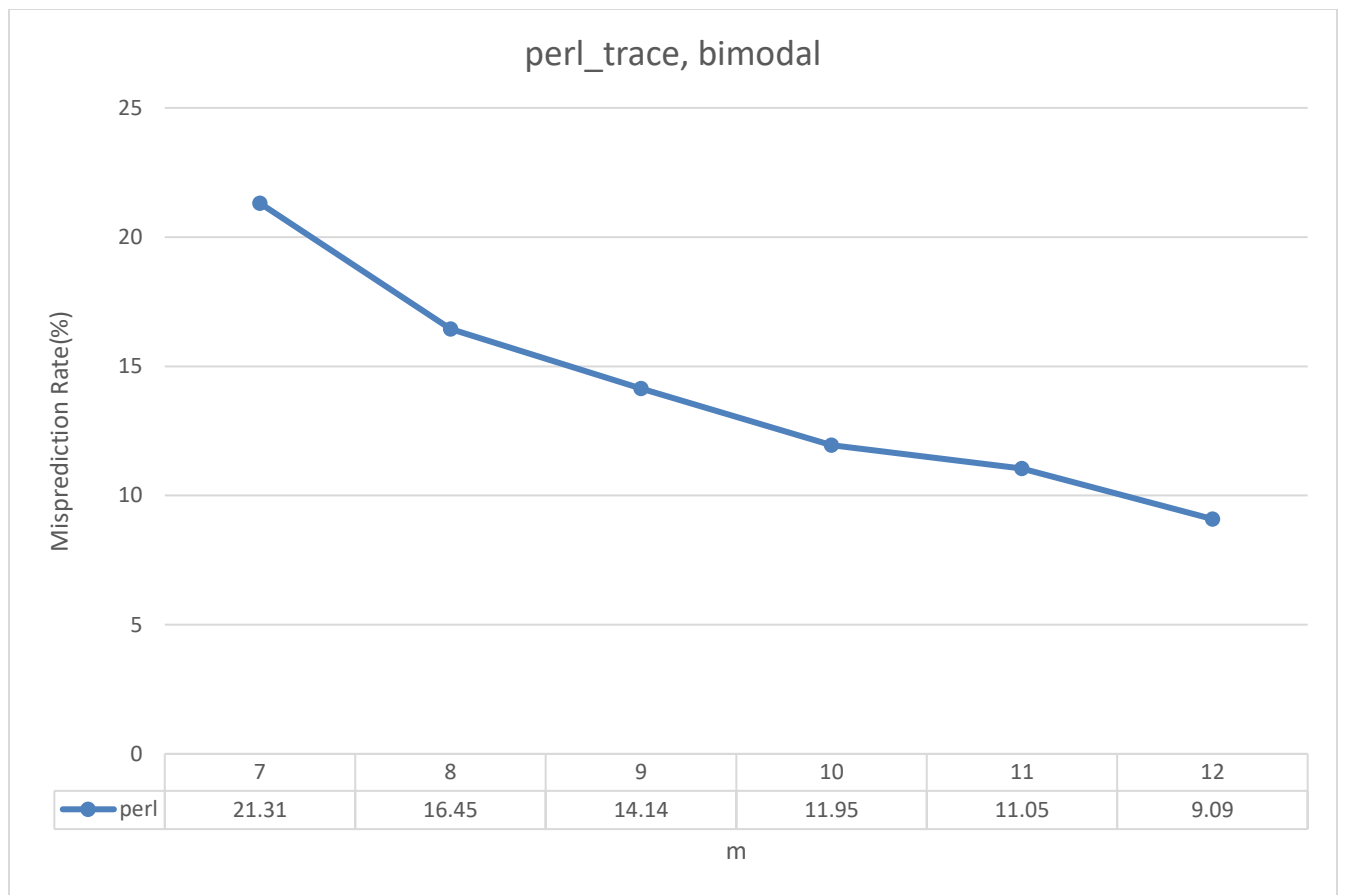
We see that in all the benchmarks, as the number of bits (m) used to index the program counter (PC) increases, the misprediction rate also drops substantially. This is because the number of entries in the prediction table depends on the value of m . Therefore, if m bits are used to index the PC, the number of entries in the prediction table equals two raised to the power of m . More the number of entries, more the number of Smith Counters which correspond to a particular index. Therefore, this index will hold more information about the instruction than an index with lesser number of bits. This will lead to a decrease in interference which will in turn lead to decrease in mispredictions.

For example, a small value of m will result in the same index for multiple branches. Hence, the counter will be updated whenever one of the branches arrive at the predictor which is not always undesirable but can lead to more number of mispredictions. Whereas with a larger index, the counter will be updated fewer number of times because now lesser number of branches correspond to a particular branch. Increasing the size of the index bits will increase the number of Smith counter which will add to the hardware cost. Therefore, an optimal number of bits should be used for the index after examining the misprediction rates and the hardware cost.

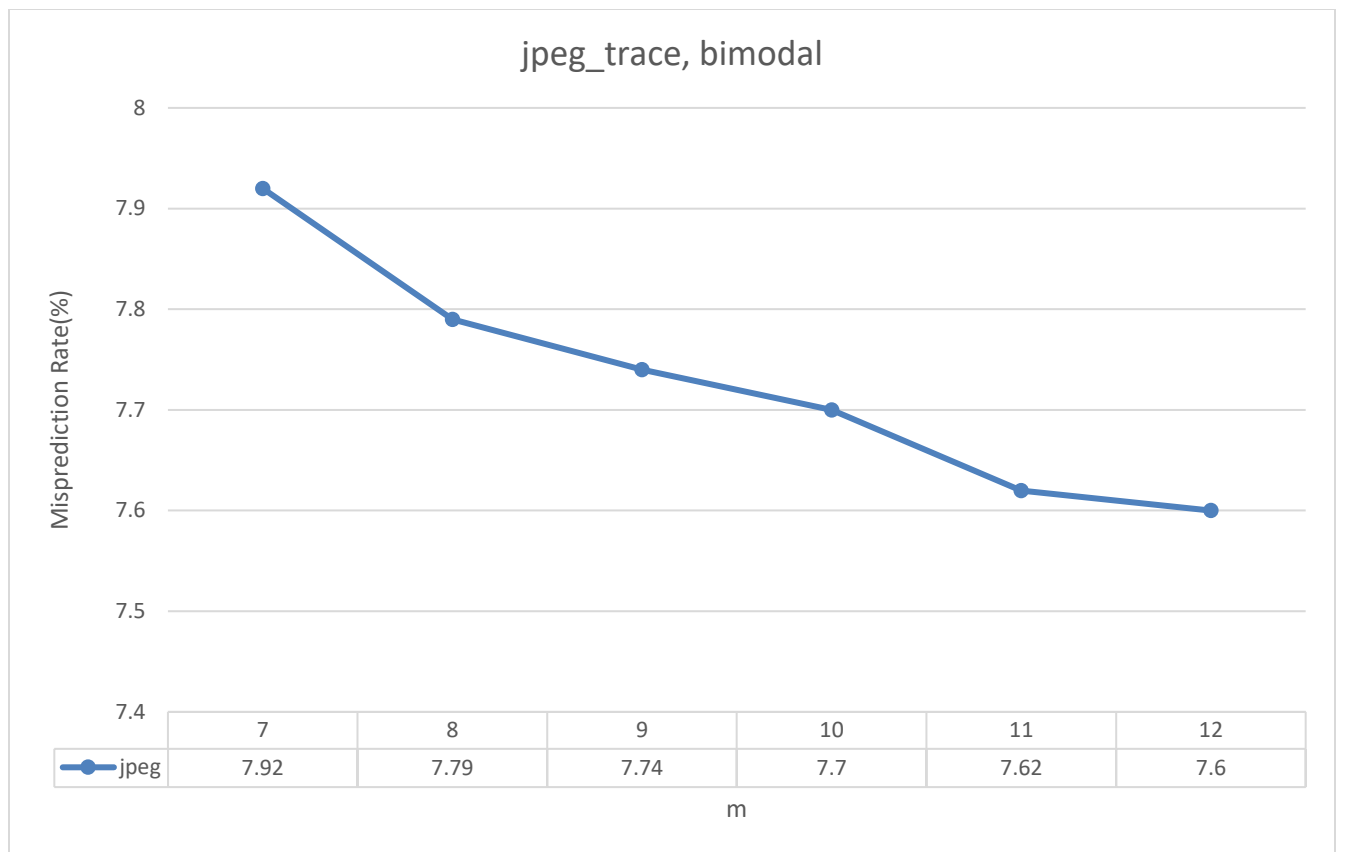
A bimodal branch predictor gives about 70-80% correct predictions. This needs to be a bigger number considering the pipeline bubbles that result from these predictions. In the case where predictions are not made, the pipeline does not suffer from any bubbles if the branch is not taken. However, when a wrong prediction is made the pipeline will suffer from a bubble no matter the branch was taken or not. The bimodal branch predictor can also give a very high misprediction rate in case of some repeating branch patterns. To correct this error, a Branch History Register (BHR) is used which tracks the history of incoming branches. Such a BHR is implemented in the gshare branch predictor in later sections. The graphs plotted from the simulation of a bimodal predictor for three different benchmarks are given below.



The gcc_trace has the highest number of mispredictions for a particular m out of all the benchmarks. A misprediction rate of 26.65 is considerable and this benchmark needs to use a different predictor or a bigger m to get satisfactory prediction results. Also, the lowest misprediction rate of this benchmark is greater than the highest misprediction rate of jpeg_trace. This shows that the choice of predictor also depends on the incoming instructions and this point must be taken into account when choosing a particular architecture.



The perl_trace has the highest rate of decrease in misprediction rate till $m=10$. After that the graph suddenly seems to flatten out for $m=11$ and then again decrease considerably for $m=12$. This shows that there are by increasing the index size till $m=10$, the branches are split up and a single counter corresponds to fewer number of branches. But after going from $m=10$ to $m=11$, very few branches are split up. It depends which bit of the branches changes after a common set of bits. The m corresponding to that bit will show a maximum dip in misprediction rate.



The jpeg_trace has the lowest number of misprediction rate for a particular m among all the traces. This suggests that same branches are repeatedly either being taken or not taken. The curve takes a steeper dip after m=10. This shows that the extra bit in the m=11 simulation separates two frequently occurring branches and thus increases the accuracy of prediction. The same thing is seen happening for the transition from m=7 to m=8.

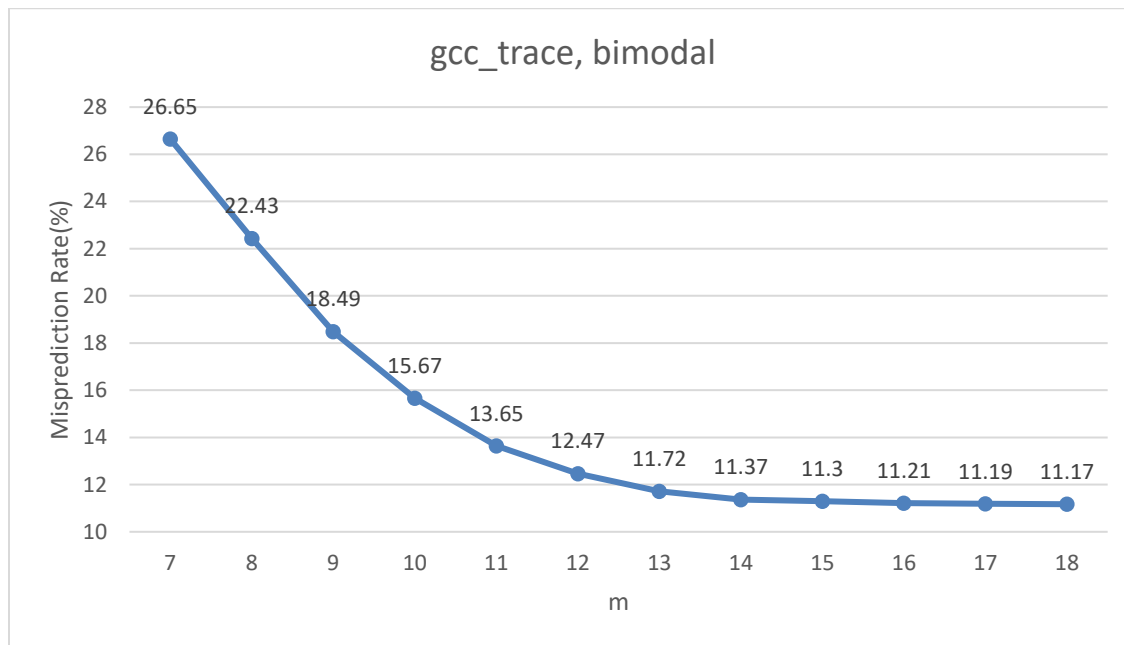
Finding optimal architecture (minimize misprediction rate and predictor cost):

Maximum Budget: 16kB

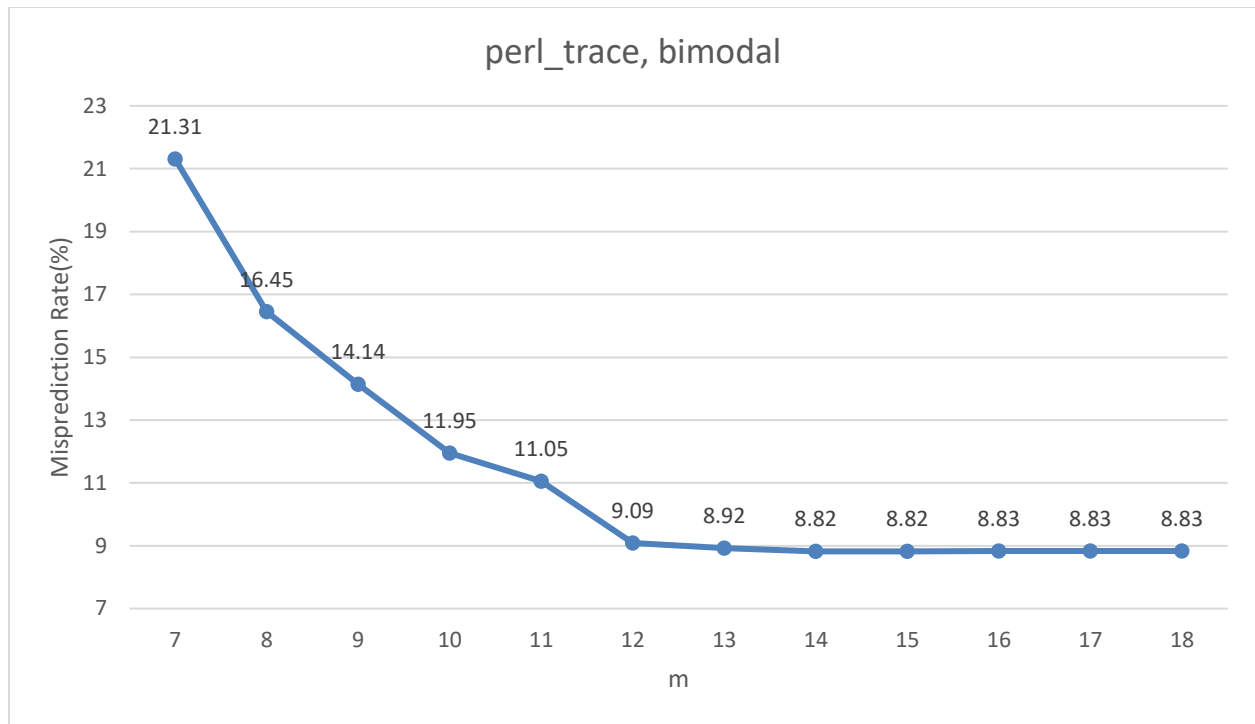
Therefore, the maximum value of m is

$$2 * 2^m = 16 * 2^{10} * 2^3$$

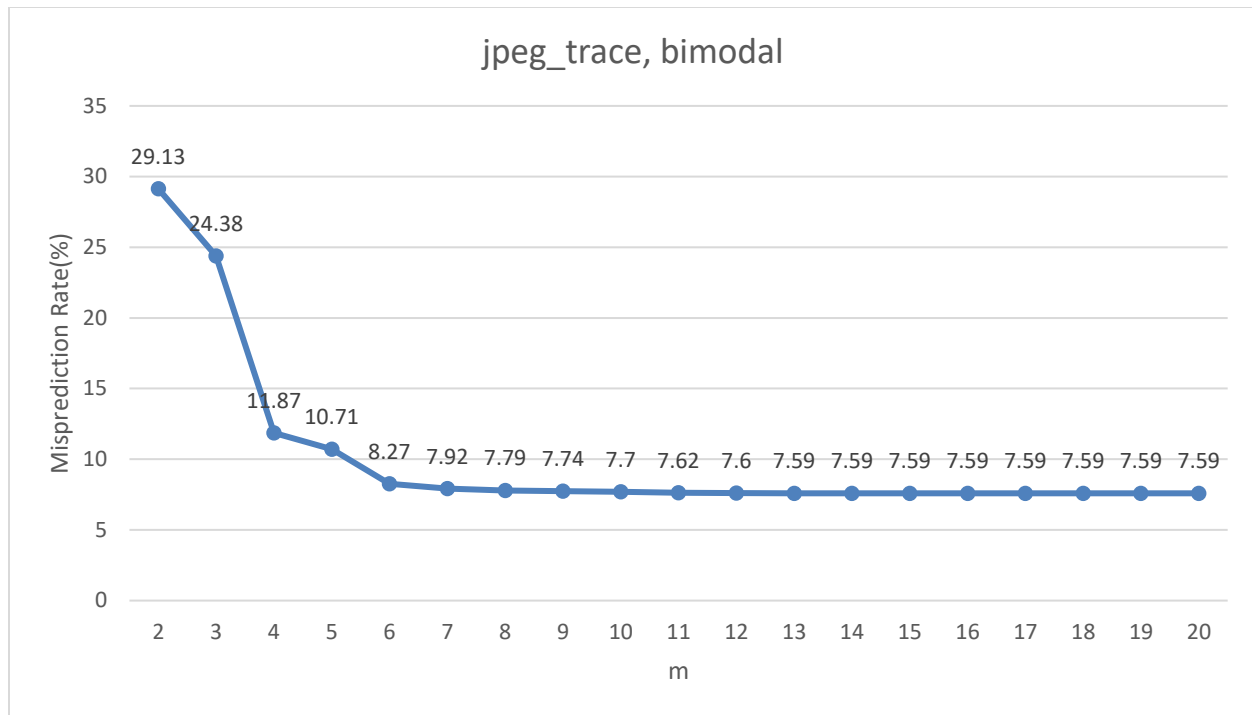
m=16



The graph becomes almost flat at m=13 which has a misprediction rate of 11.72%. Therefore, we can get a minimal misprediction rate, **11.72%** at a reasonable storage of **2kB** with **m=13**.



The slope of the graph becomes almost zero at $m=12$ where the misprediction rate falls from 11.05 at $m=11$ to 9.09 at $m=12$ and then decreases very slowly. Therefore, a good design would be to use a **$m=12$** architecture having a storage of **500B** and misprediction rate of **9.09%**.



A good design would be where $m=6$ since the values of misprediction rates become almost constant after this value. Therefore, it would be of no use to increase the hardware cost after this. Hence, the best architecture for this benchmark would be **$m=6$** with a storage of **16B** with a misprediction rate of **8.27%**.

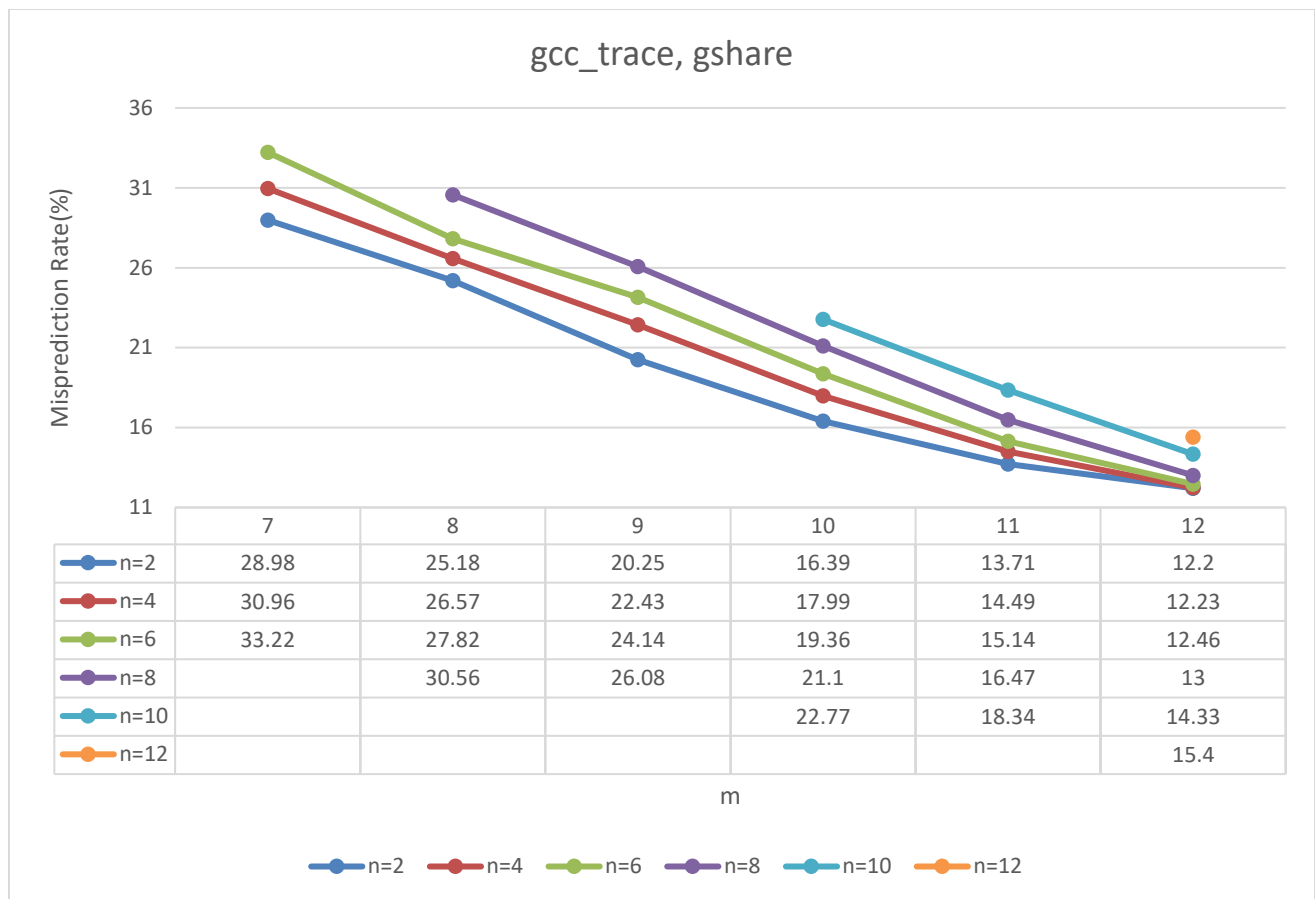
Gshare Predictor:

The difference between a gshare and a bimodal predictor is that a gshare predictor uses a Branch History Register (BHR) which is a multibit shift register. The branch outcome is shifted into this register so that it can be used to better predict the branch outcome if there is a pattern which the branches follow. The bits from the BHR are then hashed with the PC bits to form an index to the prediction table. This PC and BHR combination can also make sure that each combination goes to a different entry in the table. Therefore, the BHR makes sure that if the branches follow a pattern, the pattern is recorded and can be used to predict the upcoming branches unlike bimodal where the only evidence used to predict the branch was its previous few outcomes.

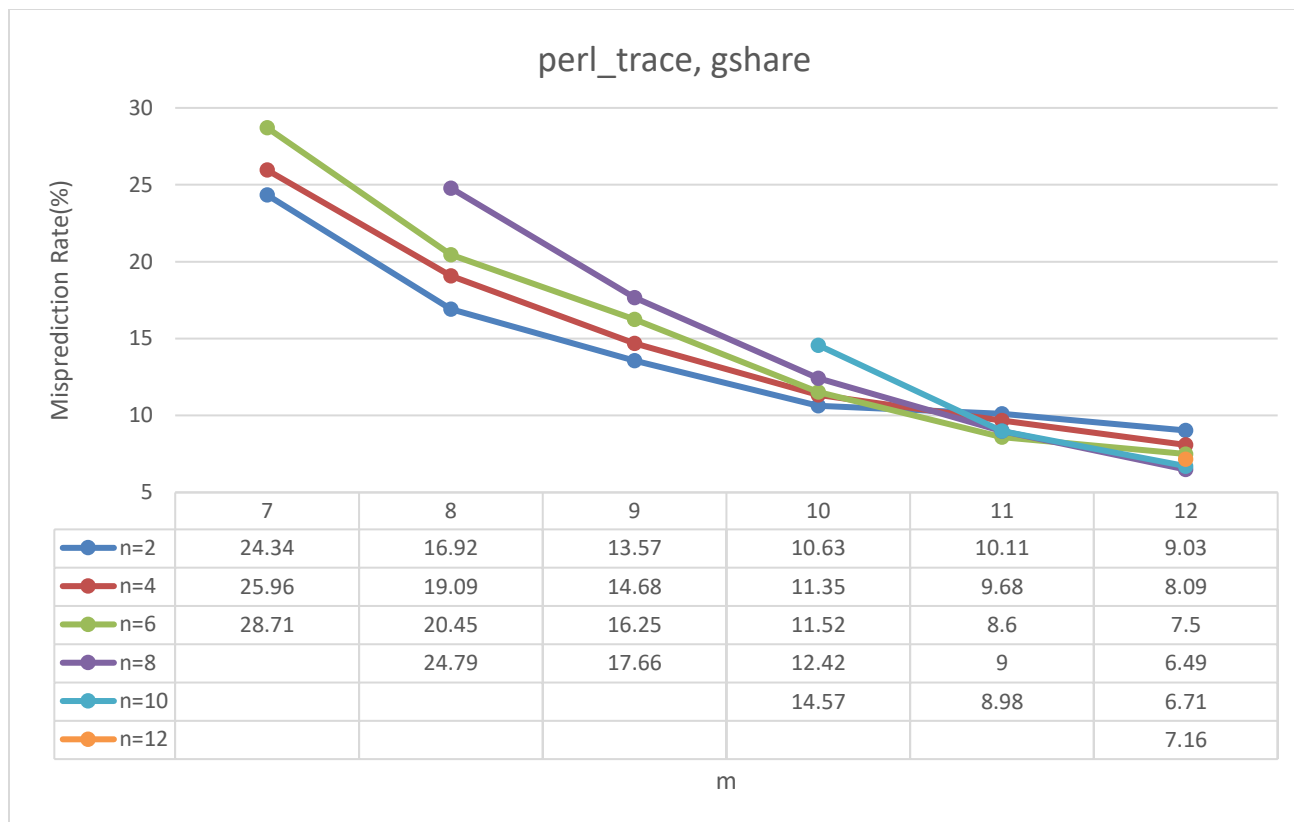
For example, if a single branch keeps flipping between being taken and not being taken, a bimodal predictor will predict based on the last few outcomes of the branch. Whereas a gshare predictor will have the pattern mapped into its BHR which will be hashed with the PC bits to give an index corresponding to the entry of either the branch being taken or not.

It is better to use bimodal if local history is more beneficial and better to use gshare if global history is more beneficial. A gshare will have more number of mispredictions if many different branches are mixed up.

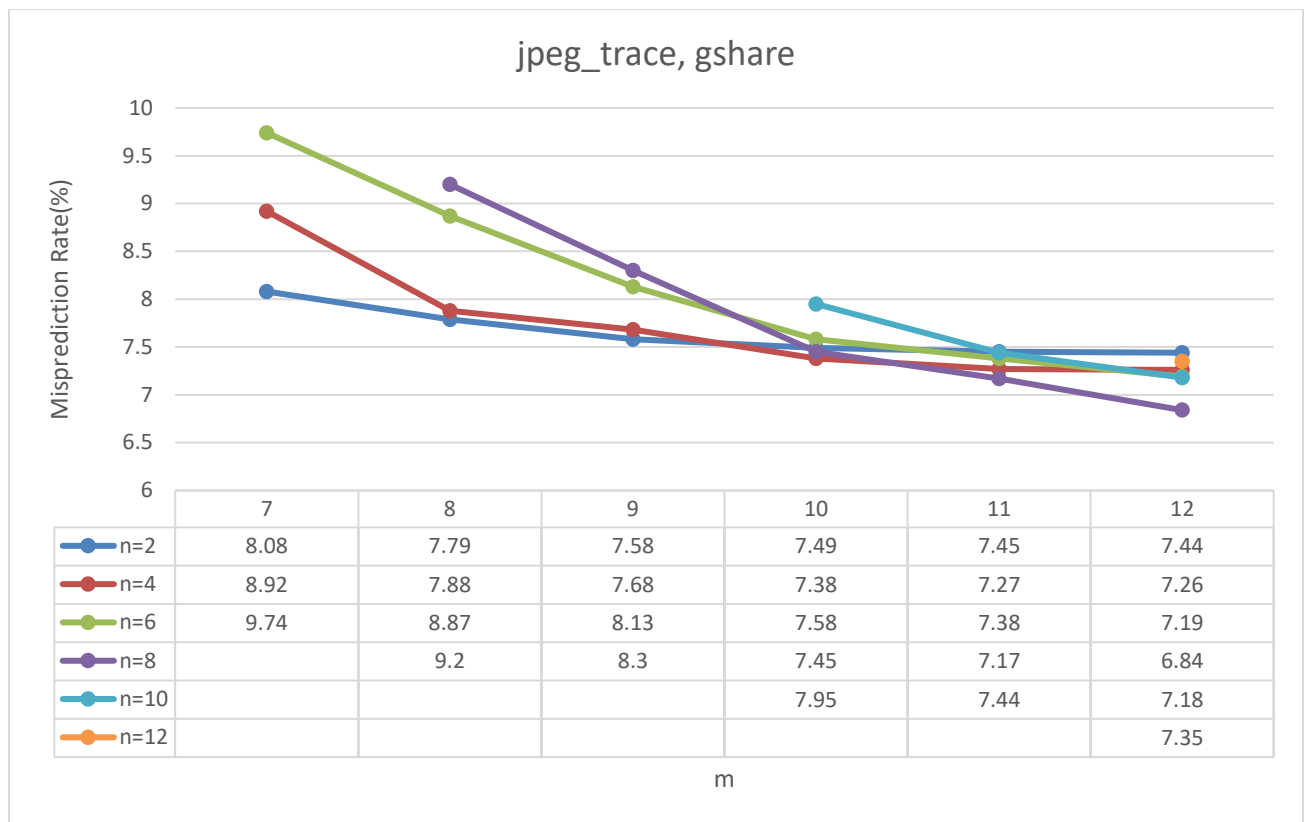
The graphs plotted from the simulations of the gshare predictor for three different traces are given below.



As the number of bits (n) used for BHR increase, the misprediction rate goes on increasing for values of m ranging from 7 to 12. It is also seen that misprediction increases at a faster rate for smaller values of m . This shows that many different are mixed up and therefore the BHR stores the history of branches which may not be useful to predict a random different branch. A small value of m also increases interference further increasing the number of mispredictions. This benchmark has the maximum misprediction rate among all the benchmarks. More the number of different branch, more the misprediction rate. Occurrence of same branch many times will decrease the misprediction rate.



The misprediction rate in this benchmark keeps increasing on increasing n for values of m from 7 to 10. For $m=11$, the value first decreases till $n=6$ and then increases to settle to an almost constant value. For $m=12$, the value decreases till $n=8$ and then drastically increases. This benchmark has an irregular variation of misprediction rate with n unlike `gcc_trace` benchmark. The lowest value of misprediction rate is observed at a configuration of $m=12$, $n=10$ for the given calculations. The highest increase in misprediction rate for a particular m is observed at $m=8$ from $n=6$ to $n=8$.



This benchmark gives the lowest values of misprediction rate among all the benchmarks. The misprediction rate decreases very slowly with increasing value of m . This benchmark shows similar results as compared to the perl benchmark where the misprediction rate increases with increasing n for $m=7,8,9$. But unlike perl benchmark, the value keeps increasing and decreasing for $m=10,11$. For $m=12$, the value decreases till $n=8$ and then increases again.

Finding optimal architecture (minimize misprediction rate and predictor cost):

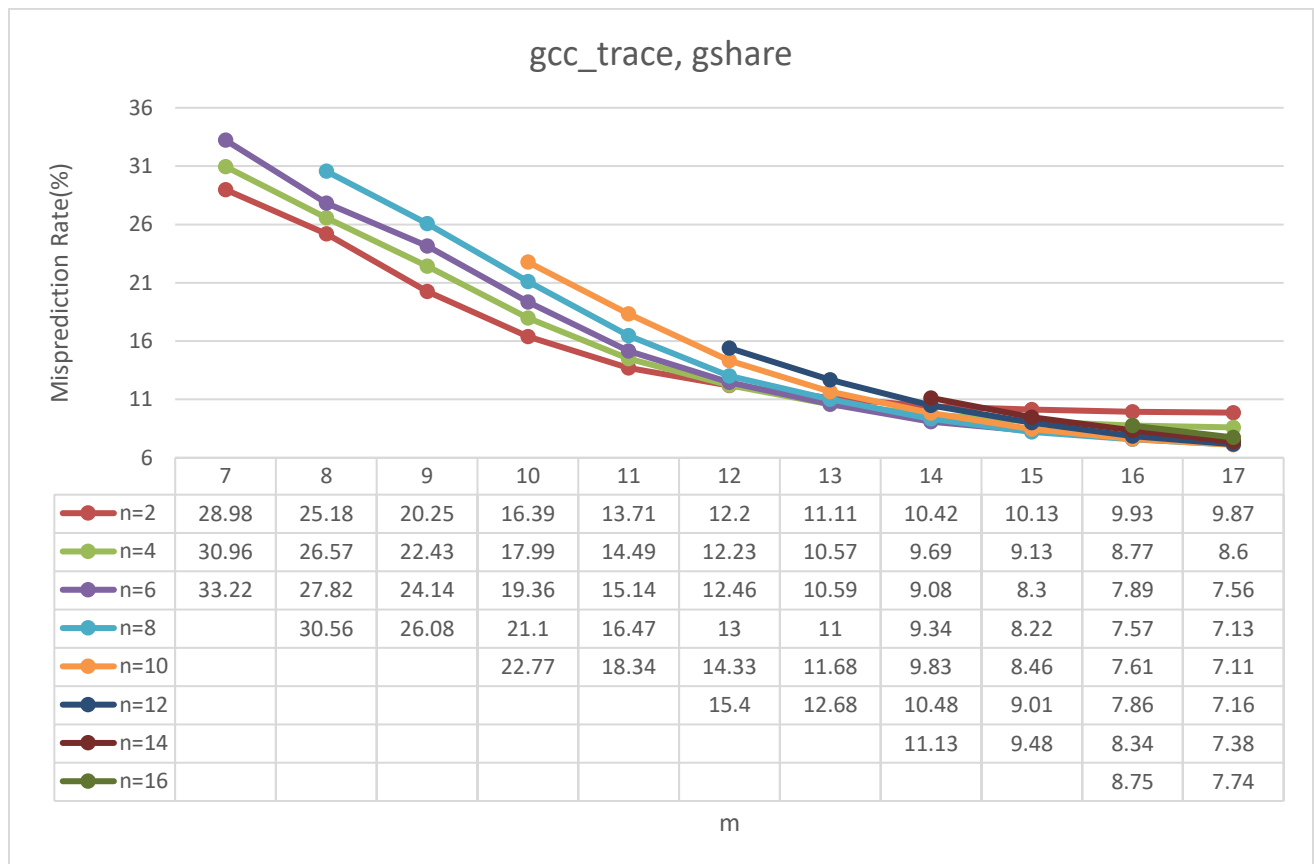
Maximum Budget: 16kB

Therefore, the maximum value of m is

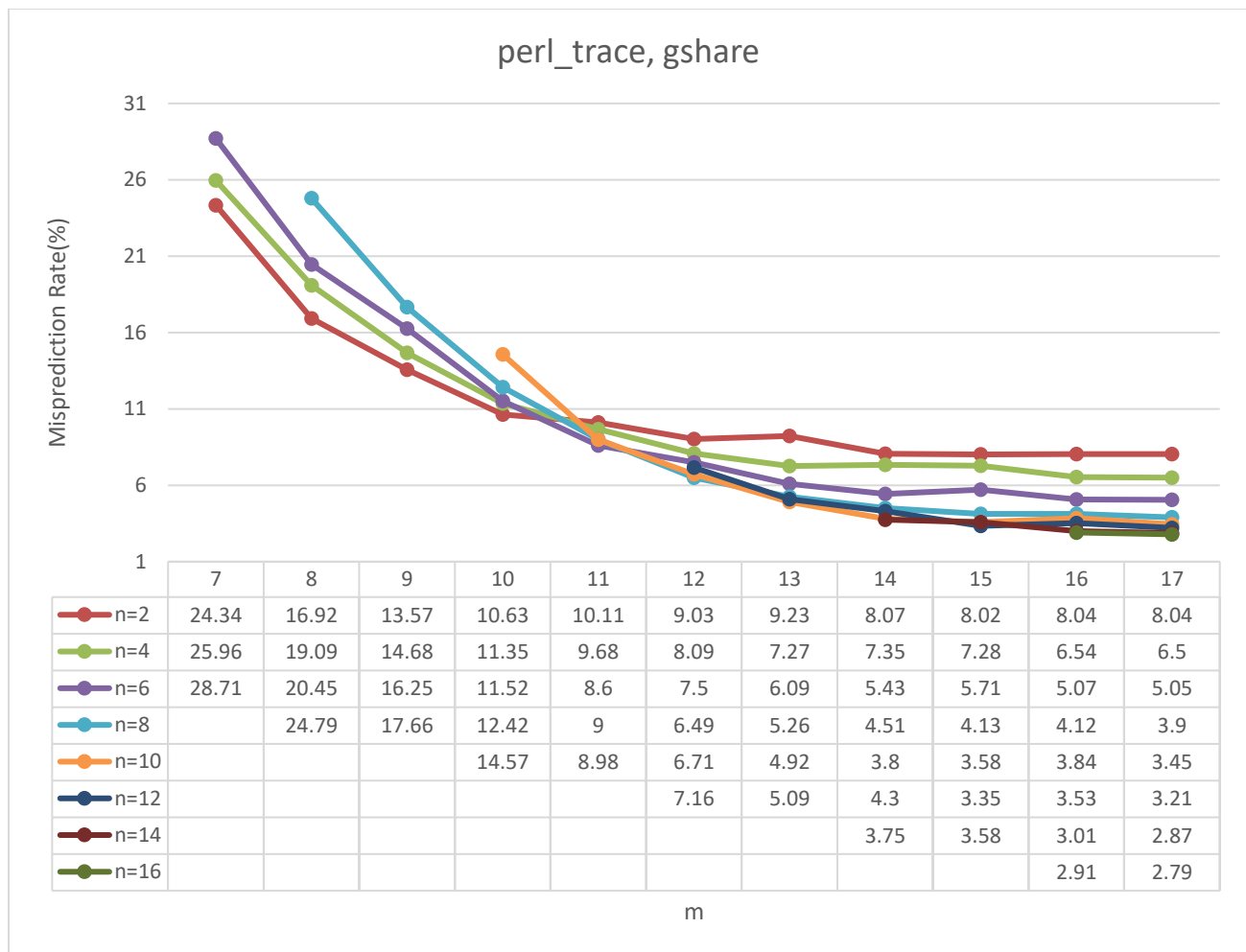
$$2 * 2^m = 16 * 2^{10} * 2^3$$

m=16

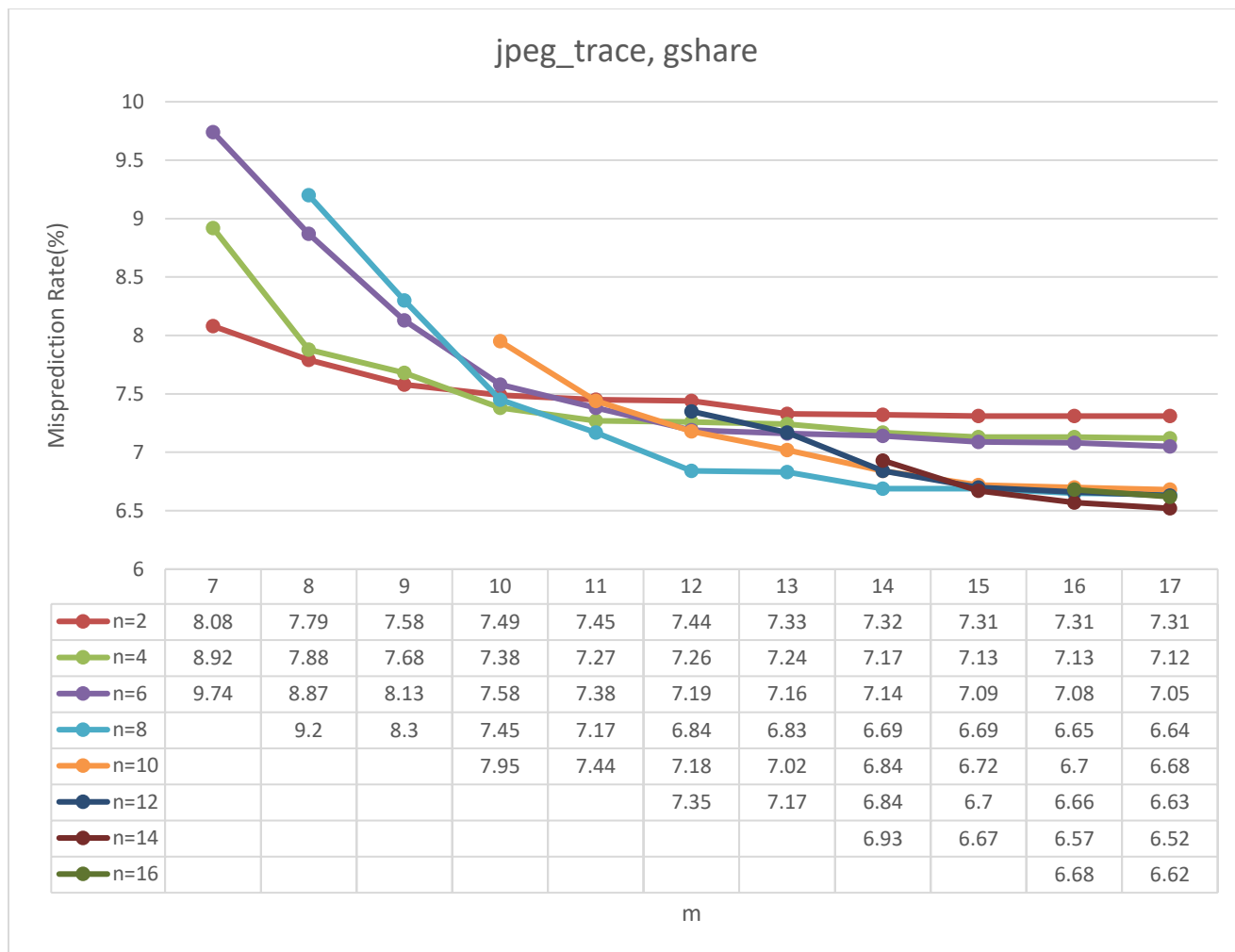
The number of bits from BHR are ignored since n is a very small number as compared to the storage of the prediction table.



The design m=16, n=8 gives a minimum value of misprediction rate. Although, if we use m=15, n=8 we get a misprediction value which is only a little greater than the one for m=16, n=8. This reduces the predictor cost. Therefore, the optimal design is **m=15, n=8** having a storage of **8kB** and misprediction rate of **8.22%**.



A good design for this benchmark would be $m=15$, $n=12$. The misprediction rate for $m=14$, $n=12$ increases by a little more than 1%. The design $m=14$, $n=10$ is also a good candidate and is better than $m=15$, $n=12$ because of the reduced predictor cost and only a small difference in misprediction rate. Therefore, the design which minimizes misprediction rate and predictor cost is **$m=14$, $n=10$** having a storage of **4kB** and misprediction rate of **3.8%**.



The graph for $n=8$ is lowest for mid-range values of m . Therefore, the design with **$m=12$, $n=8$** minimizes both the misprediction rate and the predictor cost and has a misprediction rate of **6.84%** with a storage of **1kB**.