

NC State University
Department of Electrical and Computer Engineering
ECE 463/563: Fall 2021 (Rotenberg)
Project #2: Branch Prediction

by

RAMACHANDRAN SEKANIPURAM SRIKANTHAN

NCSU Honor Pledge: "I have neither given nor received unauthorized aid on this project."

Student's electronic signature: Ramachandran Sekanipuram Srikanthan.

Course number: 563

Part 1: Bimodal Predictor

The bimodal predictor maintains a counter table to predict whether the instruction is a branch and not. If it is predicted as branch, then it predicts whether the branch should be taken or not. The predictor uses M bits of the instruction from the PC. Since M bits are being used, the bimodal predictor's counter table has 2^M elements in it. Since the lower two bits in RISC-V instruction set are 0's it can safely be ignored. Thus, we use the value from the second bit from LSB till M bits from the address traces to find the index of counter table.

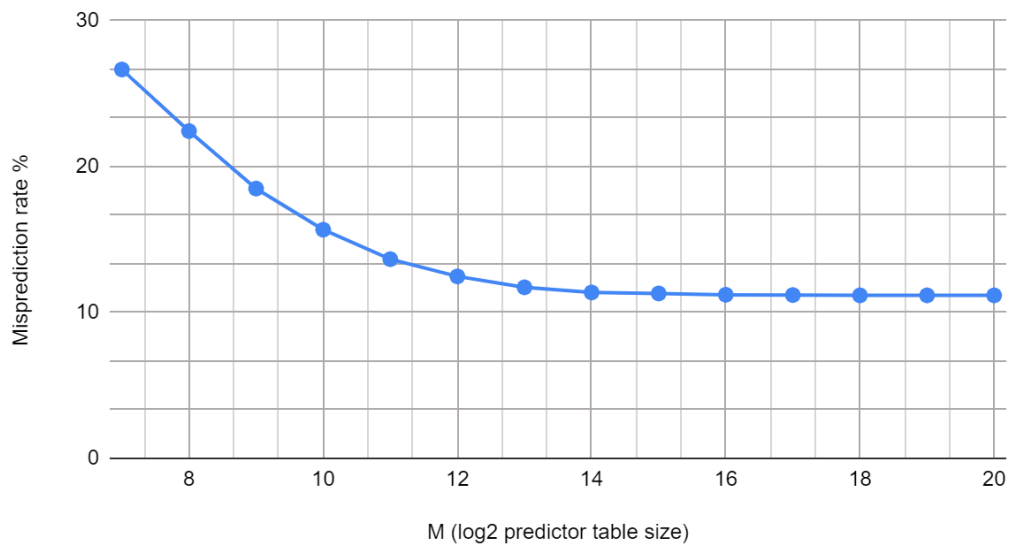
The counter table can contain value between 0 and 3. The values 0 and 3 represent strongly not taken and strongly taken respectively. The values 1 and 2 represent weakly not taken and weakly taken respectively. Initially all the elements in counter table are being set to 2 (i.e., weakly taken).

The index calculated from the previous step is used to find the prediction from the counter table. If the index value is greater than or equal to 2, it is predicted as branch taken, else as not taken. The counter table is finally updated based on the actual outcome of the branch. If the actual outcome is taken, the counter table's value is incremented till 3 as maximum. If the actual outcome is not taken, the counter table's value is decremented till 0 as minimum.

Performance of Bimodal Predictor for “gcc_trace” file

M	Misprediction rate %
7	26.65
8	22.43
9	18.49
10	15.67
11	13.65
12	12.47
13	11.72
14	11.37
15	11.3
16	11.21
17	11.19
18	11.17
19	11.17
20	11.17

Misprediction rate % vs. M (bits being used).



From the graph we can find out that, the misprediction rate decreases as we increase the value of M which directly affects the predictor table size. As M increases the number of entries in the predictor table increases and so we have a lot of space to store more information. It can also be seen that from the `gcc_trace.txt` file the number of addresses targeting the same index in the prediction table is higher for small values of M, thus increasing the misprediction rate. Thus, when M is increased, we have a few addresses targeting the same index in the prediction table thereby decreasing the misprediction rate.

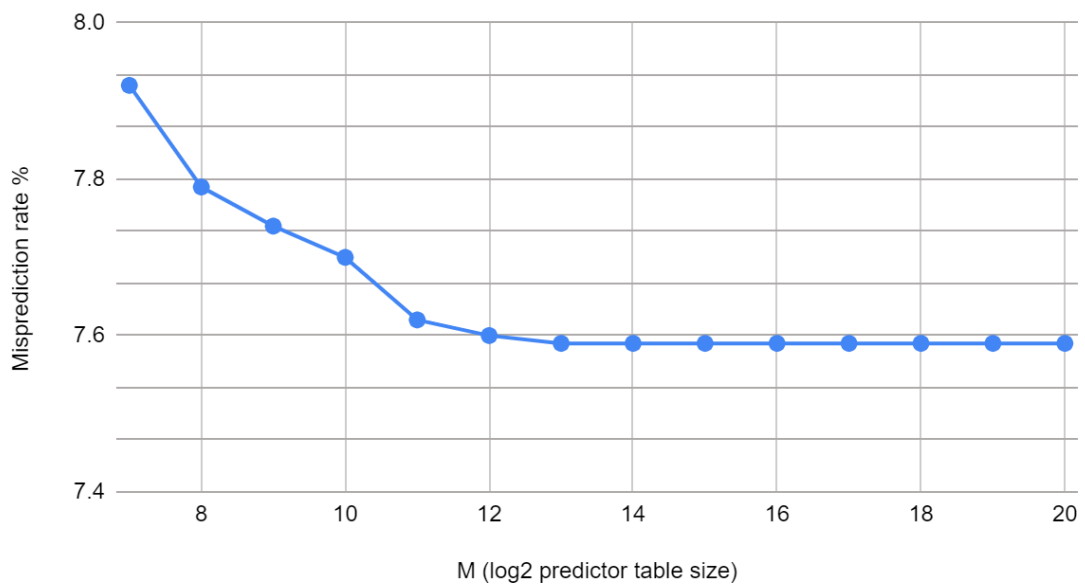
After a particular value of M ($M = 12$ in this case), the misprediction rate saturates and the rates don't diverge much as M increases. The reason for this behavior is that after a certain value of M, most of the addresses map to a unique location and so there is less collision between various addresses, thus making misprediction rate depend purely on the previously trained data causing saturation of misprediction rate.

Thus, for $M = 12$, the trade off between predictor table size and misprediction rate provides the best value.

Performance of Bimodal Predictor for “jpeg_trace” file

M	Misprediction rate %
7	7.92
8	7.79
9	7.74
10	7.7
11	7.62
12	7.6
13	7.59
14	7.59
15	7.59
16	7.59
17	7.59
18	7.59
19	7.59
20	7.59

Misprediction rate % vs. M (bits being used).

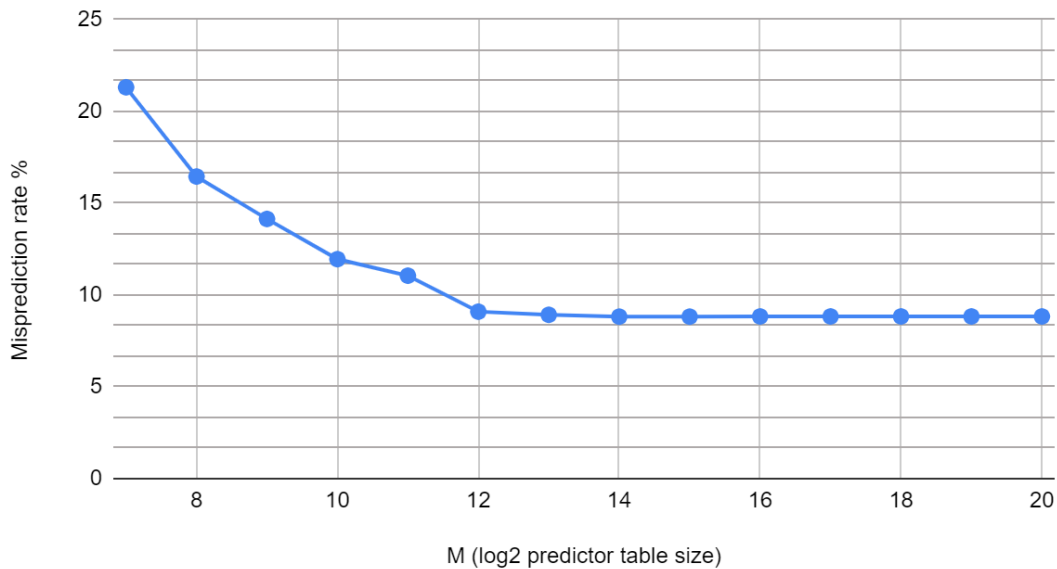


From the graph, we can see that the misprediction rate decreases as we increase the predictor table size (increasing M bits) and finally saturates after a certain value of M. We can also note from the graph that, the misprediction rate decreases steadily until it reaches a saturation value, indicating that most of addresses map to a unique index in the predictor table thereby making misprediction rate depend mainly on the previously trained data. The point of $M = 11$, gives the best value for tradeoff between the prediction table size in memory and the misprediction rate.

Performance of Bimodal Predictor for “perl_trace” file

M	Misprediction rate %
7	21.31
8	16.45
9	14.14
10	11.95
11	11.05
12	9.09
13	8.92
14	8.82
15	8.82
16	8.83
17	8.83
18	8.83
19	8.83
20	8.83

Misprediction rate % vs. M (bits being used).



The graph indicates that as we increase the value of M the misprediction rate decreases until a certain M value and becomes saturated after that. It can be seen clearly from the “perl_trace” dataset that the misprediction rate is very much high for lower values of M and a slightly increase causes a drastic increase in performance. This means that there are lot of address traces mapping to the same target index when very few bits are chosen to find out the target index. With a slight increase in M value, much of the address traces maps to a unique target index. The value of $M = 12$ gives a best value for tradeoff between memory occupancy for storing the predictor table and the misprediction rate.

Comparison among various benchmarks

From all the above benchmarks we can clearly see that the misprediction rate decreases as we increase the size of M causing a large size of predictor table in memory until a point and after that the value saturates. The point before which the misprediction rate saturates is the best point for handling the tradeoff between memory size needed to store the prediction table and misprediction rate.

We can also note that, the “gcc_trace” and “perl_trace” has a misprediction rate significantly higher than the “jpeg_trace” benchmark indicating that, the “jpeg_trace” benchmark has a lot of unique address trace, thus enabling the misprediction rate to depend mainly on previously trained data.

From the graphs, we can also note that for “gcc_trace” and “perl_trace” benchmarks the curve saturates for $M = 12$ and for “jpeg_trace” the curve saturates at $M = 11$ indicating the uniqueness in the address traces in the “jpeg_trace” benchmark.

Part 2: G-Share Predictor

The G-Share predictor also utilizes a predictor table to store and train from the previously resulted outcomes. The predictor uses M bits from the address trace and creates a table of size 2^M to store the predictions. The lower two bits of the address traces are ignored as they are 0's for RISC-V instruction set. The Predictor also utilizes a global branch history register to keep in track of branch operation in a global scope. The global history register stores the value of the actual outcomes of the branches and can remember up to N bits.

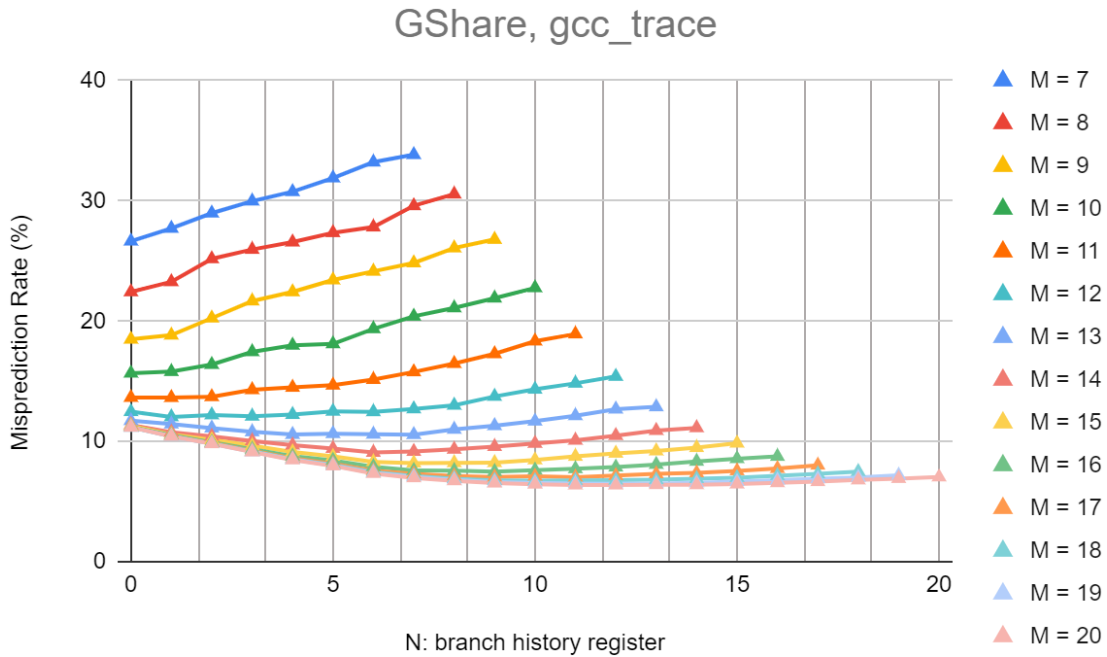
The Predictor table can be in any of the four values between 0 and 3 at any time. The values 0 and 3 represents strongly not taken and strongly taken respectively. The values 1 and 2 represents weakly not taken and weakly taken respectively. The predictor table is initialized to 2 (weakly taken). When the outcome of the corresponding address trace indicated by the index of the predictor table is taken, then the value in the table is incremented by 1 up to a maximum value of 3. When the outcome of the corresponding address trace indicated by the index of the predictor table is not taken, then the value in the table is decremented by 1 up to a minimum value of 0.

The index is calculated by using a hashing function which maps the address traces to a corresponding range of hash values which can be used to index into the predictor table. The last N bits (from MSB) of the address from the trace is XORed with the N bits of global branch history register's contents. This result is then concatenated with the remaining $M-N$ bits from the address trace to map into the index of the predictor table.

Depending on the actual outcome of the branch the global history register is also updated. The register is right shifted by 1 and the value of 1 is appended into the MSB position if the actual outcome is a branch taken, else a value of 0 is appended into the MSB.

Performance of G-Share Predictor for “gcc_trace” file

[illegible]



From the graph we can understand that, the size of global history registers also has an impact in the misprediction rate. For larger values of M , as the size of the global history register is increased there is a significant drop in the misprediction rate until a certain value. Now for a fixed value of M , we can see that as we increase the size of global history register (N), there is a significant decrease in the misprediction rate. The reason being that smaller values of N can store less global branch taken outcome. The XOR in the hashing function ensures that this previous global branch taken data is taken into consideration while calculating the index into the predictor table. The concatenation operation with the remaining $M-N$ bits of the current trace is also used in the calculation of the index into the table.

Now if we increase the value of N , the misprediction rate decreases as we have more global data taken into consideration. This also means that we have sufficient $M-N$ bits concatenated, thereby providing sufficient current address trace's data. But if the N increases beyond a value, we start to accumulate large amount of older data which may not be of significant help in the current prediction. Due to this accumulation of older data, the current prediction gets biased to all the older values causing the misprediction rate to increase. Also, $M-N$ becomes smaller as we increase N causing a much biasing to previous outcomes by giving less weight to present address trace's contents.

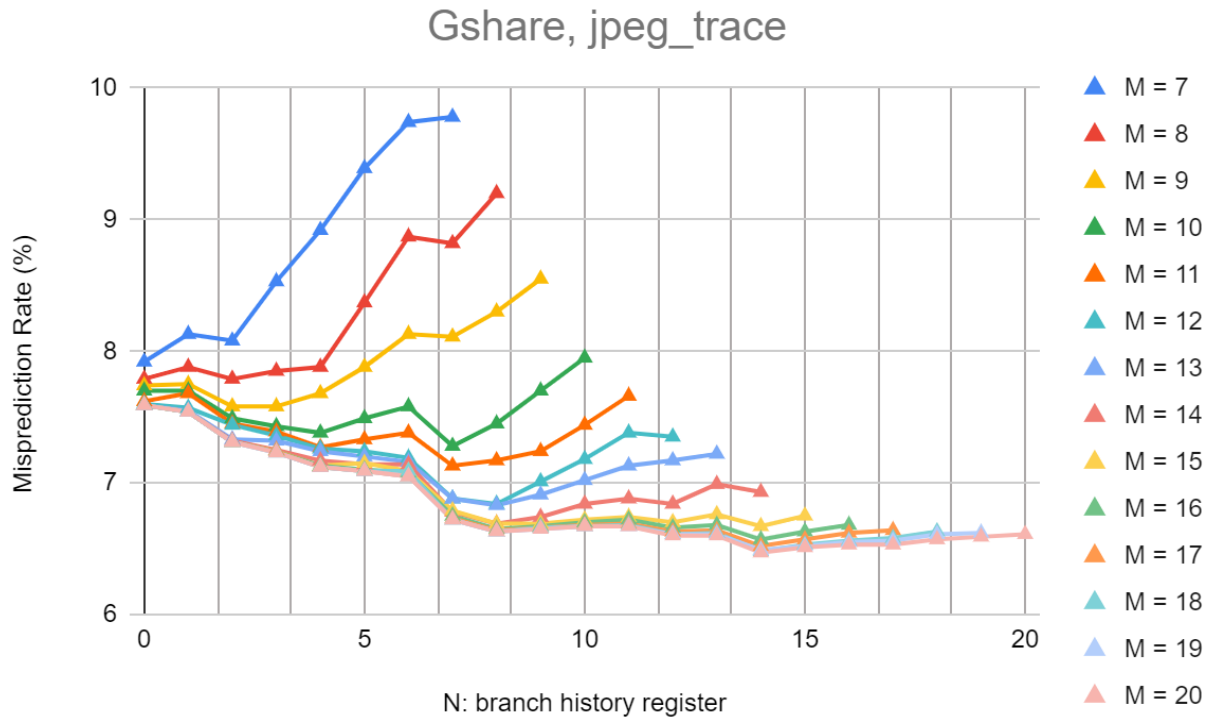
We can also see that, for smaller values of M , increasing N has a negative impact in the performance. Since for smaller values of M we may have less indexes in the table and so there may be much of collisions into the same index. In that case when we increase the N , much of older values of branch history gets accumulated giving less space ($M-N$ bits) to take into consideration of the present address's traces which of significantly smaller in size thereby increasing the misprediction rate.

Thus, for G-Share predictor we need a storage overhead is $N + 2^M$ for both predictor table and global history register. The values of M and N for which we have an optimal predictor table size and global history register size for which the misprediction

rate is not too much higher is for $M=14$ and $N=6$. Since above which for the same value of N we can see very slight decrease in the misprediction rate for which the storage overhead becomes exponentially high in powers of 2. Also, for large values of N we can see the misprediction rate goes higher as discussed before.

Performance of G-Share Predictor for “jpeg_trace” file

[illegible]



The graph indicates that misprediction rate decreases for increasing N in large values of M. The misprediction rate goes much lower for larger M values but at the cost of very large memory overhead for storing predictor table. The M and N values for which we get better performance in terms of misprediction rate with having the lesser memory overhead is for M = 12 and N = 7.

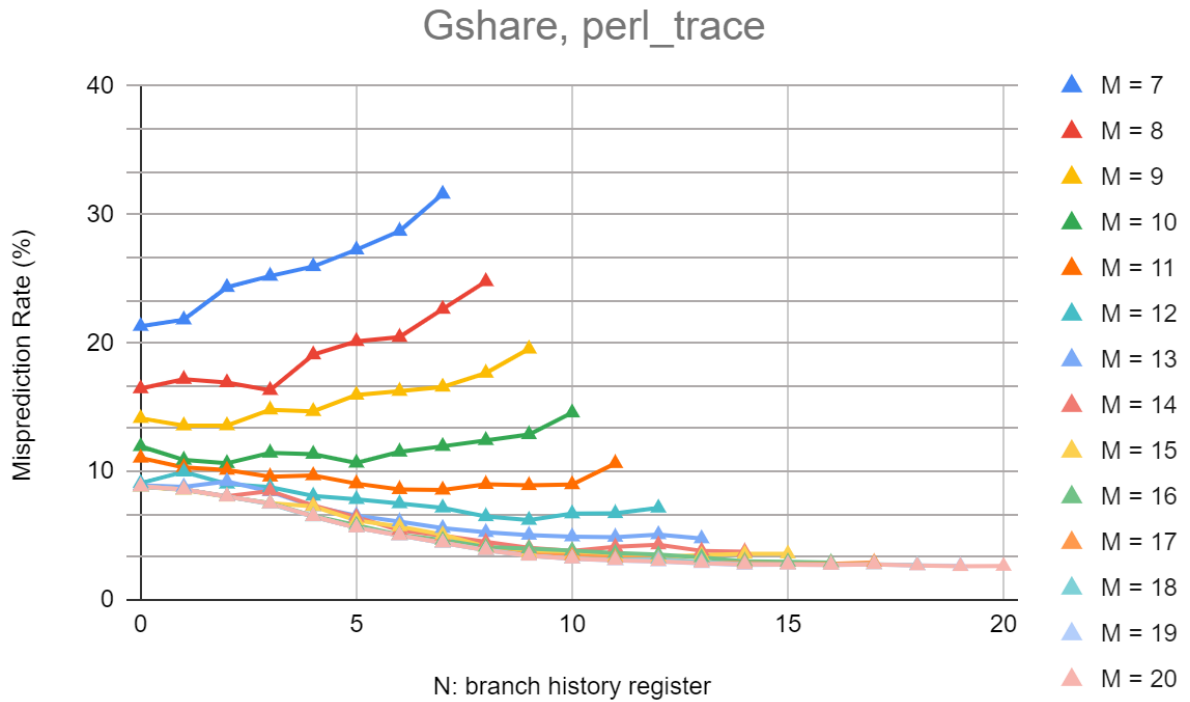
For a particular value of M (Large M value) as N increases, we can have large amount of previous global branch data up to a point and after that with increase in N the accumulation of global history data becomes higher and reduces the usage of current address trace's data thereby increasing misprediction rate. For lower values of M, increasing N causes still lowering of current address's traces values which is already small. Also, the size of predictor table is very small causing lot of mispredictions.

Here we can also see that the misprediction rate for bimodal predictor (n = 0) is lesser when compared to g-share as we increase N values for smaller m values. Thus, it is better to use bimodal predictor for small M values and G-share for large M values.

Now for large values of M, for a particular N, as we increase the value of M, we can only see a slight decrease in misprediction rate at the cost of large increase in memory overhead. Even if we increase N for large values of M, we can see curve nearly overlaps and there is very less change in the misprediction rate for M values from 12 to 20, so its better to have M = 12 as best value without compromising the tradeoff.

Performance of G-Share Predictor for “perl_trace” file

[illegible]



From the graph, we can see that for large values of M as we increase the value of N , much of global branch history is stored and is being used to predict the next outcome which thereby decreases the misprediction rate.

The best value of M and N for having an appreciable predictor table size (appreciable memory overhead) without affect too much of misprediction rate is for $M = 14$ and $N = 10$. We can further increase the value of M but the misprediction rate changes slightly but the cost of memory overhead becomes very much higher.

For a fixed value of M as we increase the N size, the history values of branches are given much weightage than the present address's traces causing increase in misprediction rate. Also, for lower values of M , the size of prediction table is small causing increase in misprediction rate.

Also, as we increase N above a point, for a particular M value, the misprediction rate increases causing degradation in predictor performance.

Comparison among various benchmarks

We can identify that for “gcc_trace” and “perl_trace” we have a value of $M = 14$ without biasing towards any side in the tradeoff between memory overhead and misprediction rate. For “jpeg_trace” we have a value of $M = 12$ indicating that this benchmark has lot of address traces mapping to unique target locations in the predictor table.

We can also note that for “jpeg_trace”, the bimodal predictor is not able to acquire the misprediction rate of below 7 even for $M = 20$, but G-share acquires misprediction rate of below 7 for $M = 12$ keeping $N = 7$. This shows that N (global history register size) has an impact in misprediction ratio. In the case of bimodal, as the

curve saturates the misprediction rate mainly depends on previously trained data for the same addresses, but here in G-share we can see that, it depends on global branch history along with previously trained data causing the misprediction rate to further reduce which cannot be obtained in case of bimodal predictor.

But for the same “jpeg_trace” benchmark, for lower values of M , bimodal predictor gives the least misprediction rate which cannot be obtained by G-Share predictor for any N as consideration of the branch taken history further negatively impacts the prediction for small M values.

Thus, having a bimodal predictor for lower M values and having a G-share predictor for higher M values with sufficient N provides better results without taking too much of memory overhead and without increasing the misprediction rate.