

User Manual of 3D Integration-Device Thermal Modeling MATLAB Framework v1.0 integrated with NeuroSim

Developers: Ankit Kaul, Yang Zhang, and Madison Manley

PI: Dr. Muhannad Bakir, Georgia Institute of Technology

December 2, 2022

1 CONTENTS

2	Introduction	1
2.1	Thermal Analysis	2
2.2	RRAM Retention.....	2
2.3	CIM Inference Accuracy Estimation.....	2
3	System Requirements (Linux).....	2
4	Installation	3
5	How to Run.....	3
6	References	6

2 INTRODUCTION

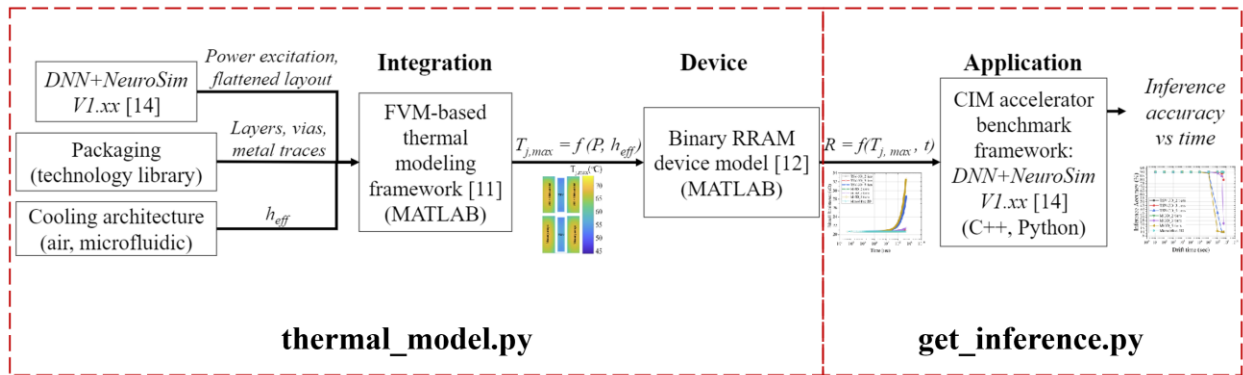


Fig. 1. Device/integration towards application-level 3D CIM reliability evaluation methodology. (Note: $T_{j,max}$ = Memory tier maximum junction temperature ($^{\circ}\text{C}$); P =Total package power (W); h_{eff} =Effective heat transfer coefficient of heat sink ($\text{W}/\text{m}^2 \cdot ^{\circ}\text{C}$); R =RRAM device resistance (Ω), t =time (sec).

This thermal modeling framework, which is developed in MATLAB and wrapped by python scripts, is used to quantify the impact of 3D integration design parameters on CIM inference accuracy [1]. This flow combines a finite volume method (FVM)-based thermal analysis

framework with a measurement-calibrated binary RRAM retention model and a CIM inference accuracy estimation framework, to perform a device-integration-application-level reliability evaluation. Details of each flow component are described as follows:

2.1 THERMAL ANALYSIS

A modified version of the FVM-based thermal modeling framework described in [2] was used to model the 3D heterogeneous integration structures and perform steady state thermal simulations. Inputs for this step include a flattened layout of the modeled CIM IP (such as memory array, ADCs, etc.), power excitation maps for each active tier/die based on the flattened layout, and a description of the die stack-up, i.e. the bulk material, interconnects, and dielectrics along with their thermal properties (thermal conductivity, specific heat capacity, etc.). For validation, the memory tier maximum junction temperatures $T_{j,max}$ from our thermal models were compared with finite-element ANSYS Mechanical APDL models (ver. 2021 R1). The maximum deviation in $T_{j,max}$ and $T_{j,min}$ between our models and ANSYS was 3.3%.

2.2 RRAM RETENTION

The memory tier $T_{j,max}$, as a function of integration design parameters such as architecture, number of tiers, input power, boundary conditions, etc., is an input to the second part of this flow that consists of a measurement-calibrated HfO_2 -based binary RRAM device model [3]. This analytical model is used to estimate device resistance variation over time, which is converted to a resistance drift ratio. The drift ratio is calculated as $\frac{R_{10}-R_{on}}{R_{on}}$, where R_{on} is the device low resistance state (LRS) resistance and R_{10} is the device resistance at 10 years operating at a specific junction temperature.

2.3 CIM INFERENCE ACCURACY ESTIMATION

Device retention change in RRAM-based CIM inference or training accelerators can correspond to a change in locally stored DNN weights/inputs/output (for weight-, input- or output-stationary dataflows [4], respectively), thus affecting the accuracy of an inference operation. The device retention drift ratio obtained from the previous step is used to adjust the drift coefficient of a retention model within DNN+NeuroSimV1.xx (a popular framework to benchmark CIM accelerators [5]), to estimate the variation of CIM inference accuracy.

3 SYSTEM REQUIREMENTS (LINUX)

The tool is expected to run in Linux with the required system dependencies installed. These include MATLAB, the matlab.engine python library, and dependencies needed for the DNN+NeuroSim Framework V1.3. We have tested the compatibility of the tool with the Linux environment Ubuntu **version** using MATLAB v2022b with python v3.8.15 and DNN+NeuroSim Framework V1.3.

The matlab.engine library is dependent on the MATLAB and python version. Please make sure that the MATLAB and python versions are compatible here or else the matlab.engine library will not be supported.

4 INSTALLATION

Step 1: Clone repo from Github

Git clone link

Step 2: Clone DNN+NeuroSim Framework V1.3 into the same directory as the Thermal Modeling software.

```
git clone https://github.com/neurosim/DNN_NeuroSim_V1.3.git
```

Step 3: Download the MATLAB engine

- Directions can be found [here](#).
- Make sure that the MATLAB and python version are compatible [here](#).

5 HOW TO RUN

1. Run thermal_model.py

After setting up the parameters, the user will run the thermal model. In this file, the user can choose to do a design space exploration. All dimensions are in micrometers. An explanation of the areas the user can explore:

0 - Single run: One simulation with the current design.

1 – chip thickness: Changes the chip thickness of all dies.

2 – chip size: Changes the chip size in the x and y dimension for all dies.

3 – TSV Diameter: Changes the TSV diameter for all dies with TSVs.

4 – TSV Pitch: Changes the chip pitch in the x and y dimension for all dies with TSVs.

5 – Bump Diameter: Changes the bump diameter for all dies.

6 – Bump Pitch: Changes the bump pitch in the x and y dimension for all dies.

Example run: `python thermal_model.py 1 200 400`

- Run the thermal model for die thickness at 200um and 400um

The program will print out the results for each simulation, which will include the $T_{j,min}$ and $T_{j,max}$ for each chip and will calculate the drift coefficient based on the device model. The screenshot below shows an example output. A **run_summary.txt** file is created that saves the last design space exploration the user has run.

```

***CALLING MATLAB ENGINE***
-----Single Run-----

Die size: 9.690e-03 m
Die area: 9.390e-05 m2
Die bulk thickness: 3.000e-05 m
TSV diameter: 1.000e-06 m
TSV pitch (x and y): 2.000e-06 m
Bump diameter: 1.800e-05 m
Bump pitch (x and y): 3.600e-05 m
Power for Chip 1: 5.050e+00 W
Power for Chip 2: 1.997e+02 W
Total stack power: 2.048e+02 W
Building stack layers
there are 10 physical layers
there are 11 nodal layers to be solved

meshing the stack
The locations of the chip are fine
The locations between the chip and package are fine
The locations between the chip and heat spreader are fine
chip domain meshing done
package domain meshing done
heat spreader domain meshing done
Elapsed time is 0.218855 seconds.
meshing done

Vias insertion
layer 1 bump information: there are 72361 vias
layer 2 bump information: there are 72361 vias
layer 2 TSV information: there are 18900 vias

build Stiffness Matrix
heat spreader top done: 0.989679

other heat spreader layers done: 1.147312
heat spreader bot done: 4.902319
other chip layers done: 3.923111
package top layer done: 2.245290
other package layers done: 0.000319
package bot done: 0.429704
11 of 11 layers are analyzed
Elapsed time is 13.653992 seconds.

assign power map
Elapsed time is 61.169307 seconds.
total power: 204.793480
finish assigning power map

calculate temperature
Elapsed time is 54.938678 seconds.
calculation done

chip1 Max: 120.735463
chip1 Min: 94.449867
chip1 Metal Max: 124.213609
chip1 Metal Min: 94.852661
chip2 Max: 128.467728
chip2 Min: 95.257298
chip2 Metal Max: 128.495830
chip2 Metal Min: 95.500884
Package Max: 128.337733
Package Min: 95.748453
total power escaped from top surface: 203.819258
Heat spreader Max: 59.307238
Heat spreader Min: 35.361616
chip1 Min ~ Max: 94.45 ~ 120.74
chip2 Min ~ Max: 95.26 ~ 128.47
total run time: 212.67

-----Calling Binary RRAM Model-----
Drift Coefficient: 7.142e-01

```

Fig. 2. Sample output from thermal_model.py

2. Run get_inference.py

After running the thermal model, the user needs to run the DNN+NeuroSim program. If the user did not place the DNN+NeuroSim folder into the thermal modeling working directory, the user will need to go into the neurosim.sh shell and change the path to the DNN+NeuroSim folder.

```

# Path to NeuroSim
cd "DNN_NeuroSim_V1.3/Inference_pytorch"

```

Fig. 3. Change the path to neurosim in neurosim.sh shell

Once get_inference.py is called, the script will get the inference accuracy from the DNN+NeuroSim program for the last design space exploration the user has run. The console will be printed to the folder **AccuracyOutput/console/<name.txt>** where **<name.txt>** is the last option the user has run with the thermal_model.py.

Last thermal_model.py run by user: python thermal_model.py 1 200 400

- Runs the thermal model for die thickness at 200um and 400um
- **python get_inference.py**
- Text files created in **AccuracyOutput/console/**
 - **1_200.txt**
 - **1_400.txt**

The console will output information in the NeuroSim simulation. The inference accuracy will be printed to the text files found in **AccuracyOutput/console/**

An example of the console output is below:

6 REFERENCES

- [1] A. Kaul, Y. Luo, X. Peng, S. Yu, and M. S. Bakir, “Thermal reliability considerations of resistive synaptic devices for 3d cim system performance,” in 2021 IEEE International 3D Systems Integration Conference (3DIC), 2021, pp. 1–5. [Online]. Available: 10.1109/3DIC52383.2021.9687612
- [2] Y. Zhang, Y. Zhang, and M. S. Bakir, “Thermal design and constraints for heterogeneous integrated chip stacks and isolation technology using air gap and thermal bridge,” IEEE Transactions on Components, Packaging and Manufacturing Technology, vol. 4, no. 12, pp. 1914–1924, 2014. [Online]. Available: 10.1109/TCPMT.2014.2364742
- [3] P.-Y. Chen and S. Yu, “Compact modeling of rram devices and its applications in 1t1r and 1s1r array design,” IEEE Trans. Electron Devices, vol. 62, no. 12, pp. 4022–4028, 2015. [Online]. Available: 10.1109/TED.2015.2492421
- [4] Y.-H. Chen, J. Emer, and V. Sze, “Eyeriss: A spatial architecture for energy-efficient dataflow for convolutional neural networks,” in Proceedings of the 43rd International Symposium on Computer Architecture, ser. ISCA '16. IEEE Press, 2016, p. 367–379. [Online]. Available: <https://doi.org/10.1109/ISCA.2016.40>
- [5] X. Peng, S. Huang, Y. Luo, X. Sun, and S. Yu, “Dnn+neurosim: An end-to-end benchmarking framework for compute-in-memory accelerators with versatile device technologies,” in Proc. IEEE Int. Elec. Devices Meeting, 2019, pp. 32.5.1–32.5.4. [Online]. Available: 10.1109/IEDM19573.2019.8993491