

Zero Load Predictive Model

Axel Jantsch

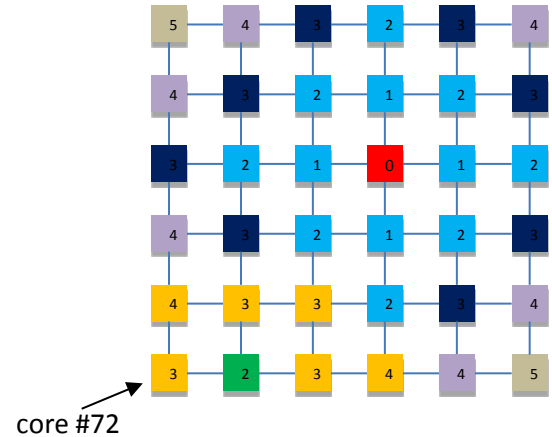
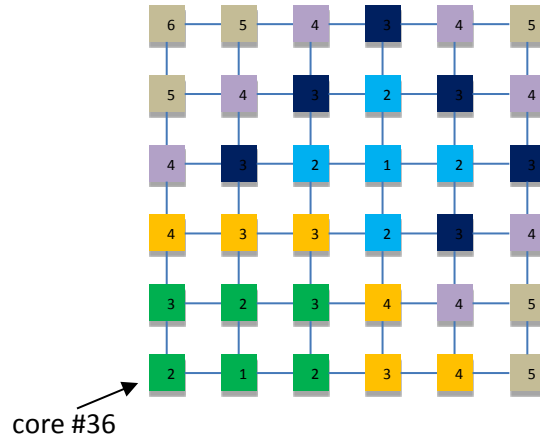
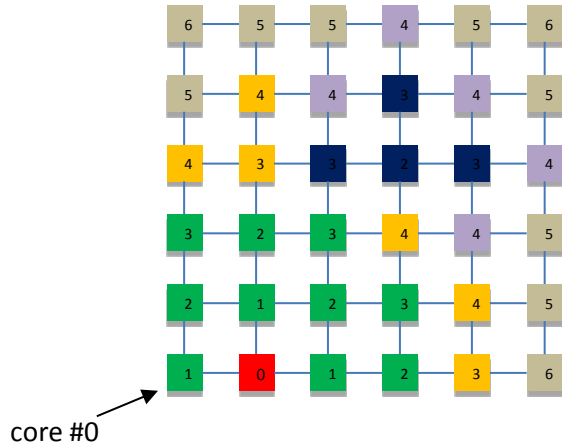
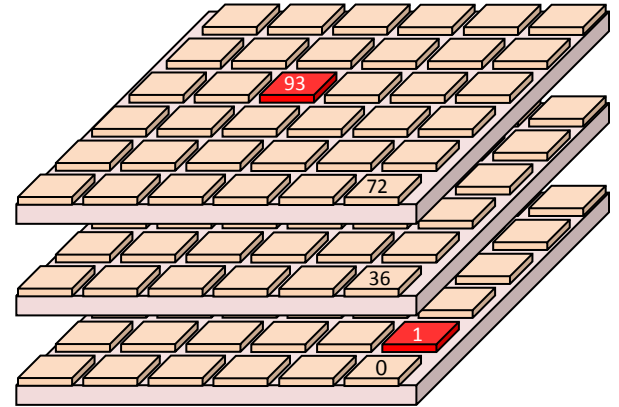
*Awet Y. Weldezion, Matt Grange, Roshan Weerasekera,
Hannu Tenhunen, Dinesh Pamunuwa,*



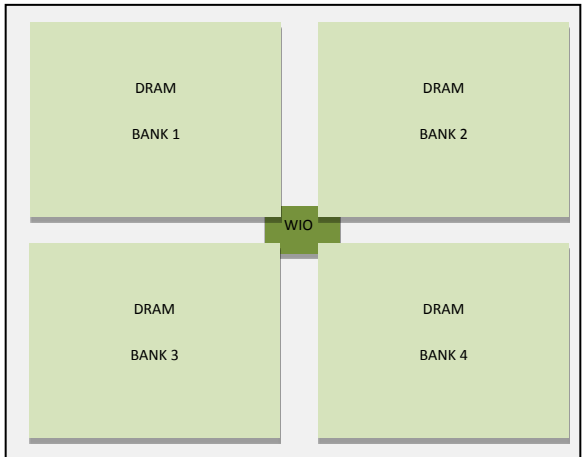
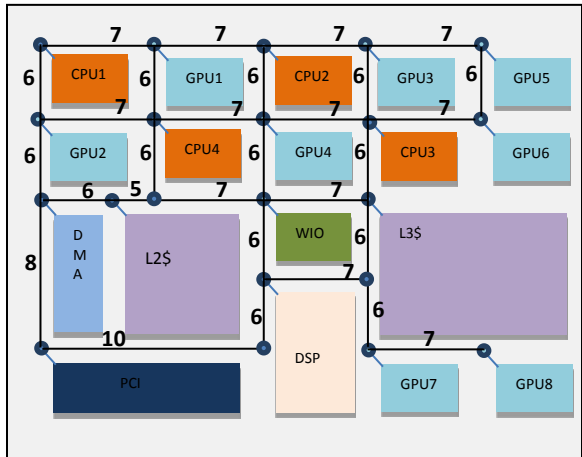
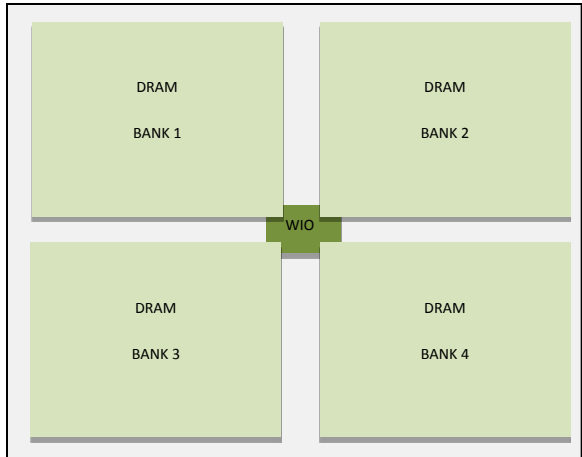
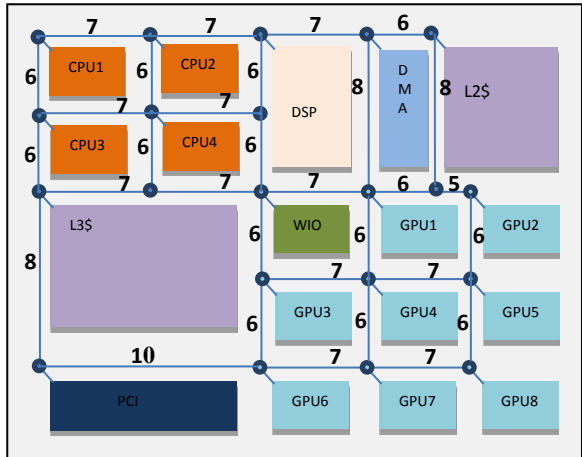
June 2013

Many Possible Configuration

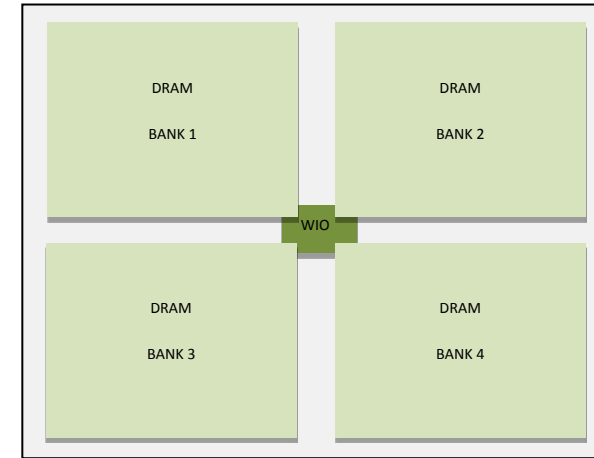
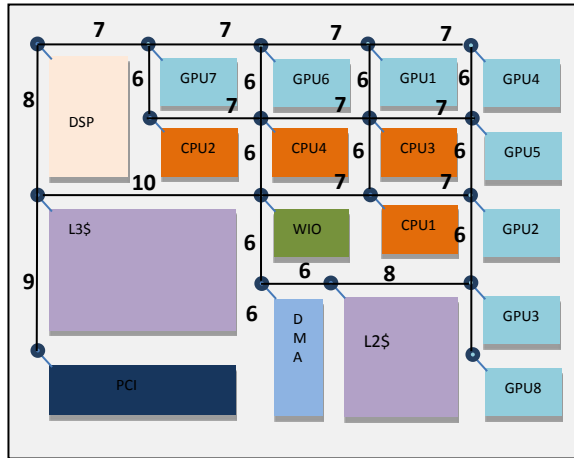
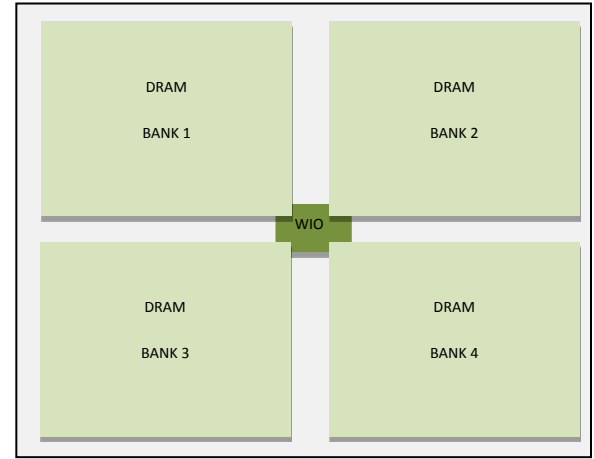
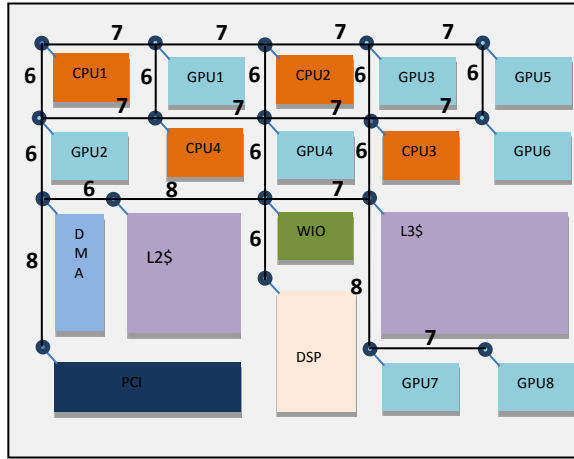
- 3 Stacked Layers with 6x6 nodes
- Two hotspots with placement constraints
- 5 types of processing units/nodes
- $\sim 5 \cdot 10^{80}$ possibilities



Heterogeneous Systems



Heterogeneous Systems



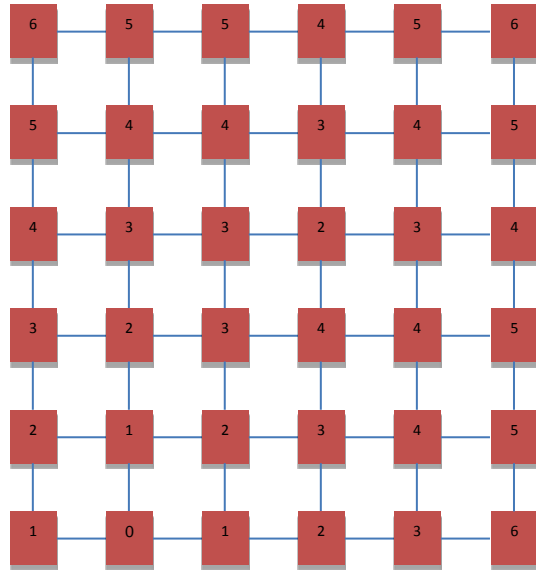
**Can we find the
best configuration
without simulation**

?

Zero Load Predictive Model

is a static model and predicts, which configuration has better performance under load.

Average distance

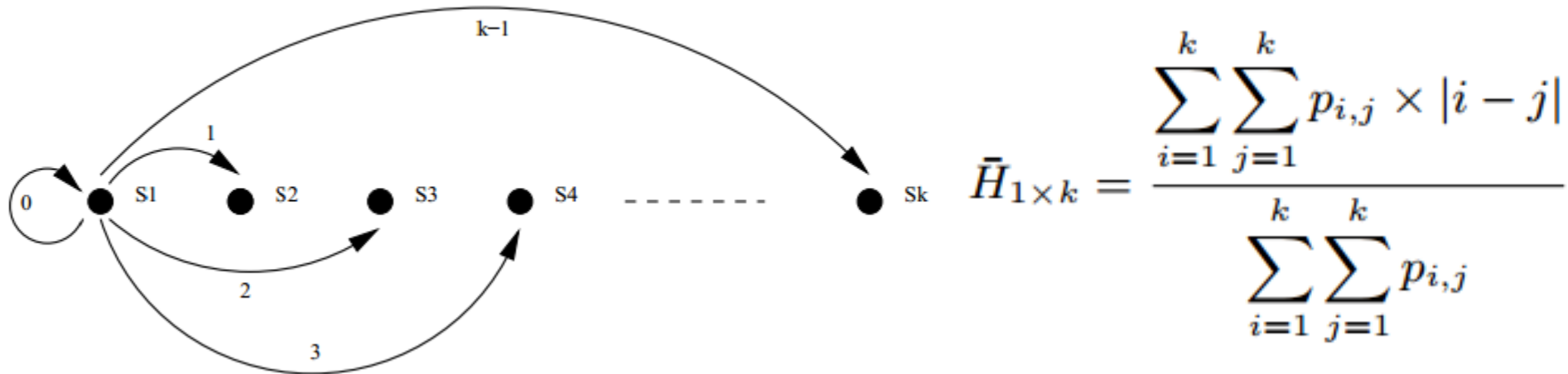


$n \times n$ Mesh for uniform random traffic:

$$\overline{H} = \frac{2}{3}n$$

Average distance

- Average distance depends on the probability, $p(i,j)$, of a packet to be sent to a destination and the actual source-destination Manhattan distance in terms of hops.



Average Distance for Uniform Random Traffic

With self-traffic

With no self-traffic:

➤ n-dimensional mesh with radix k $\overline{H} = \frac{n}{3} \left(k - \frac{1}{k} \right)$

$$\overline{H} = \frac{n}{3} k$$

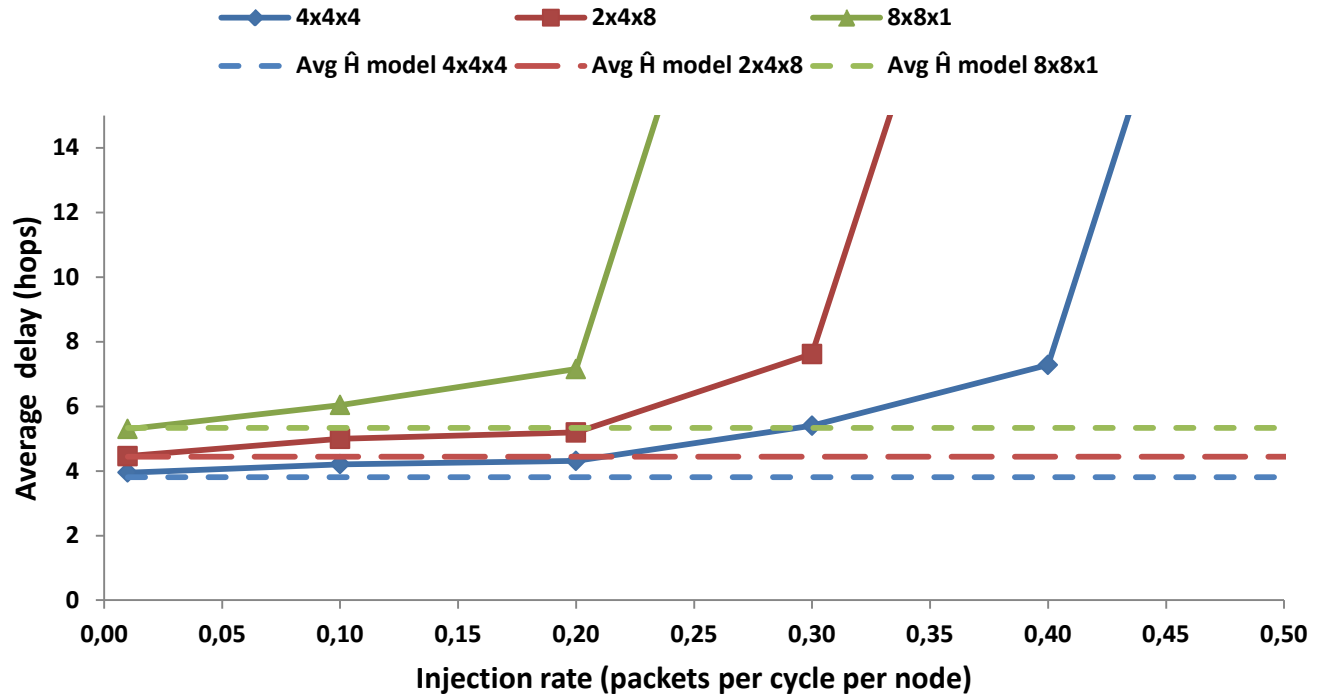
n=3 with different k in each dimension:

$$\overline{H} = \frac{1}{3} \left(k_1 - \frac{1}{k_1} + k_2 - \frac{1}{k_2} + k_3 - \frac{1}{k_3} \right)$$

$$\overline{H} = \frac{1}{3} (k_1 + k_2 + k_3)$$

Uniform Random Traffic URT

Avg. hop count for URT



Spatial Distribution

- Each source node injects a packet that is sent to at least one destination based on a pre-defined patterns such as URT, Bit-reverse, Bit-complement....
- The source address (S) is in the form of bits (b) expressed as follows

$$\hat{S} = b_m b_{m-1} b_{m-2} \dots b_3 b_2 b_1$$

Bit-Reverse

- The destination is a mirror image of the source address.
- The average distance is the average of the sum of individual source-destination distances in all directions.

$$\hat{H}_{br,xyz} = |z_x - \hat{S}_x| + |z_y - \hat{S}_y| + |z_z - \hat{S}_z|$$

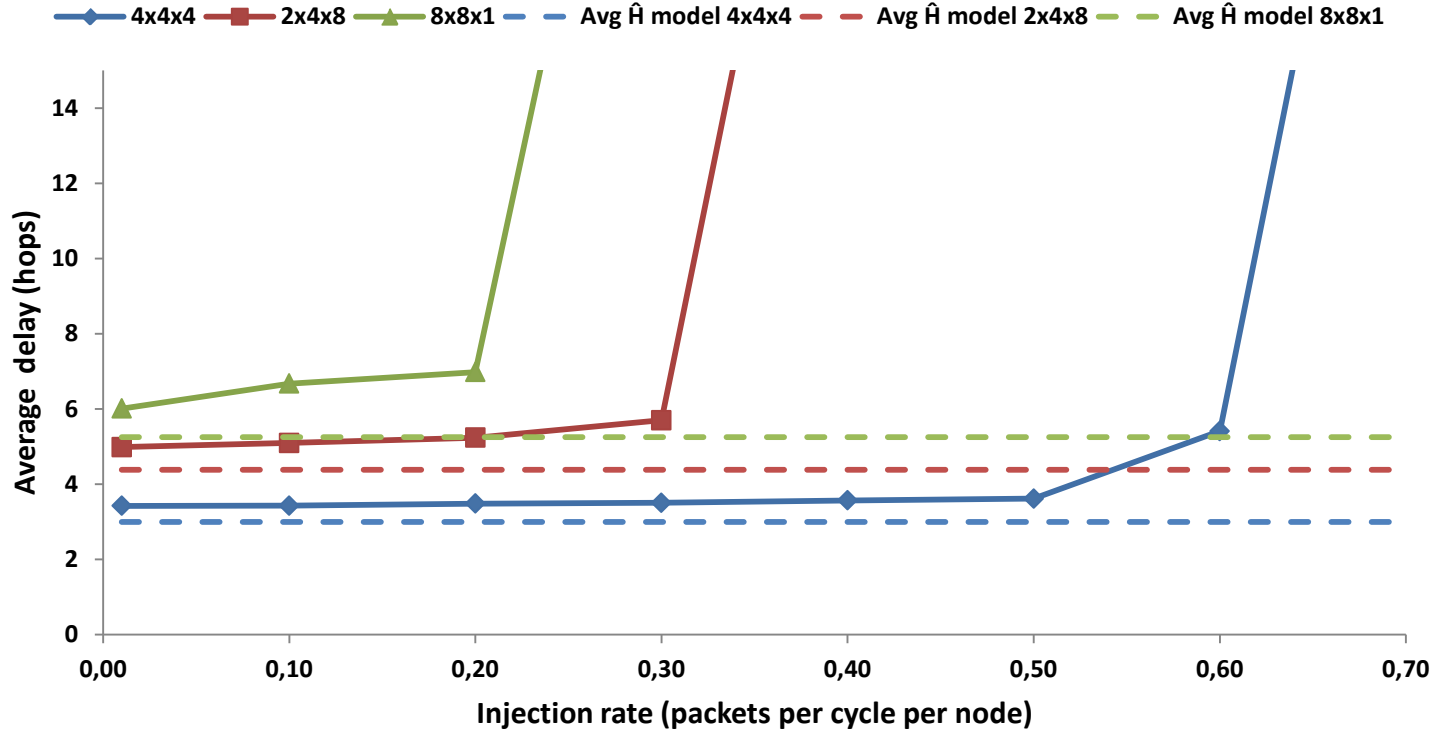
$$\hat{H}_{br}(N) = \frac{1}{N} \sum_{x \in N} \sum_{y \in N} \sum_{z \in N} \hat{H}_{br,xyz}$$

\hat{S} – Source adress

z – reversed destination

Bit-Reverse

Avg. hop count with Bit-reverse



Bit-Complement

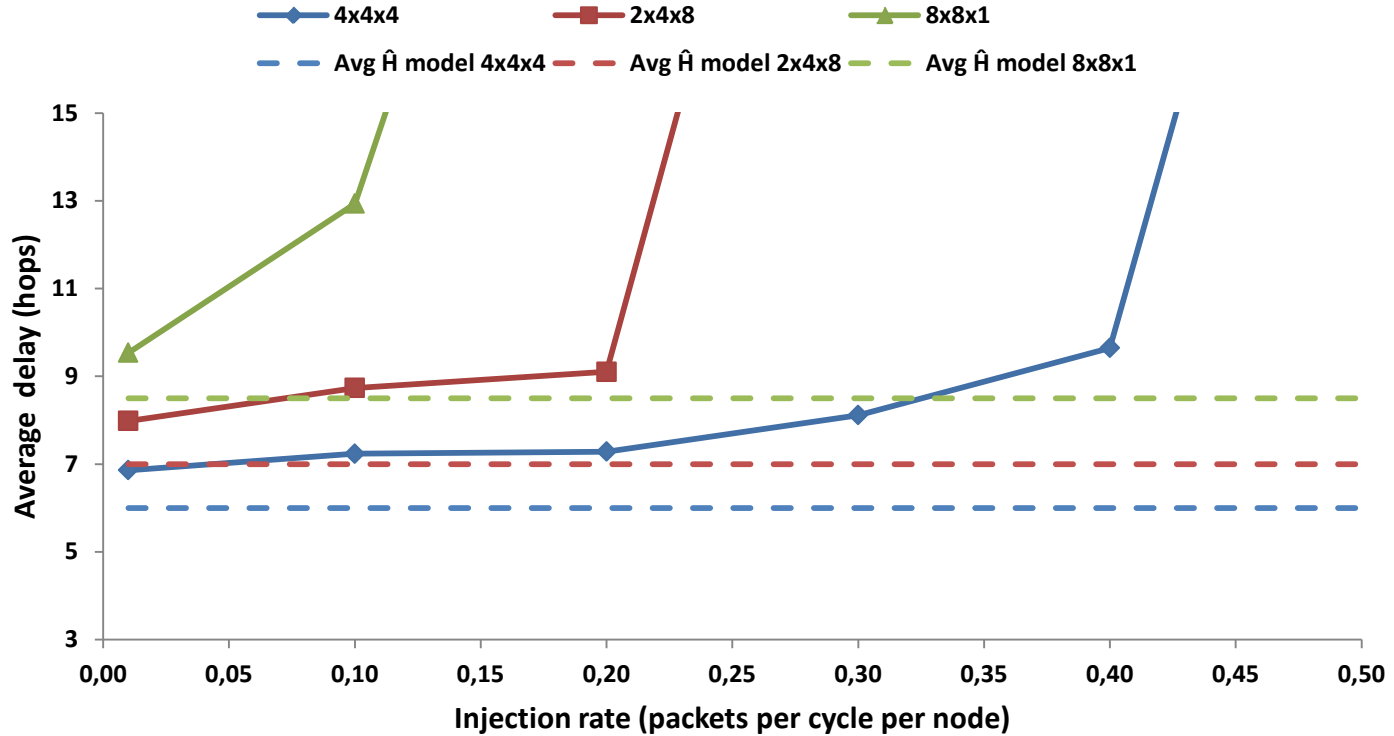
- A destination address is set as the 1's complement of each bit of the source address
- For a given 2D and 3D network the average distance is expressed as follows respectively.

$$\hat{H}_{bc,2D}(N) = \frac{x + y}{2}$$

$$\hat{H}_{bc,3D}(N) = \frac{x + y + z}{2}$$

Bit-Complement

Avg. hop count with Bit-complement



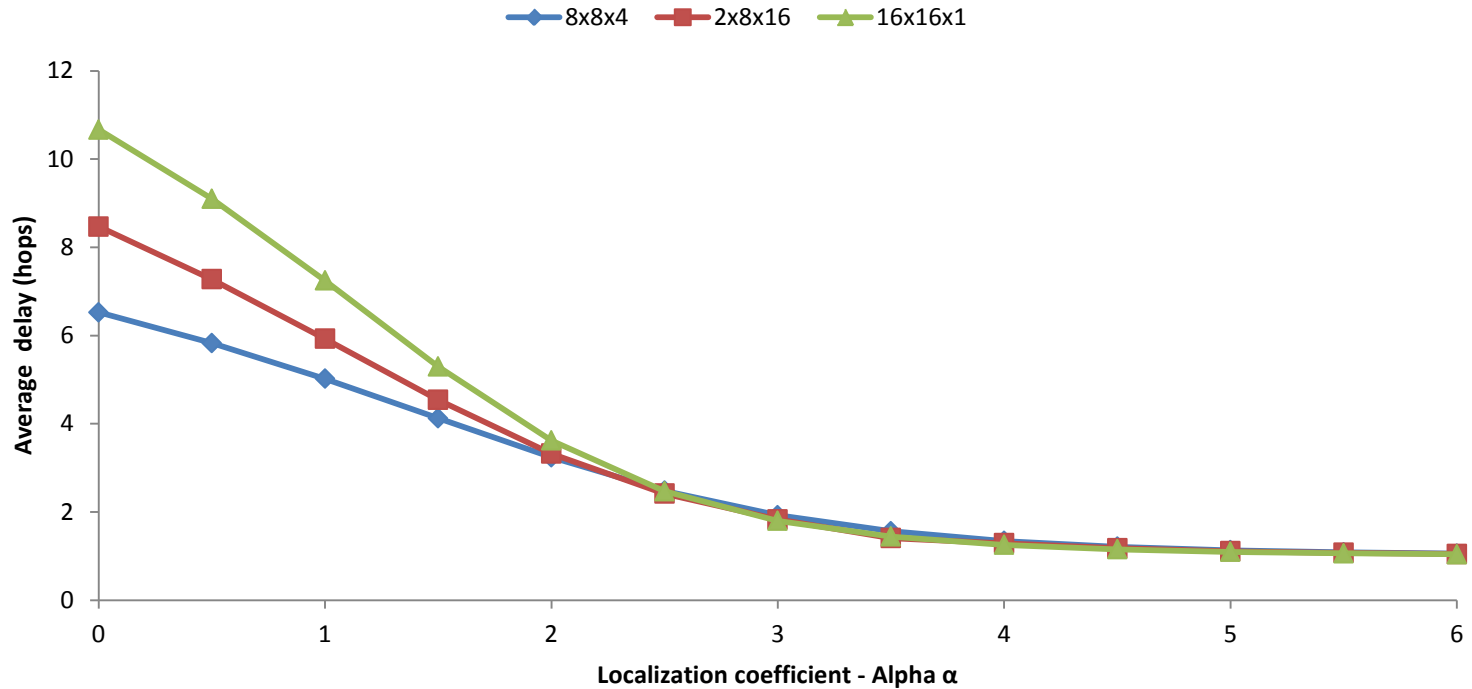
α -Model Localized Pattern

- Under a local traffic pattern, packet destinations are preferred for destinations (D) close to the source.
- The variation in level of closeness (localization) is represented with locality coefficient alpha, (α). When $\alpha = 0$ it means the level of localization is low (in this case localization does not exist)Output

$$\hat{H}_{\alpha,xyz} = \frac{1}{N-1} \sum_{\hat{S}=0}^N \sum_{\check{D}=0}^N \frac{|\hat{S} - \check{D}|^{\alpha+1}}{\sum_{\check{D}=0}^N |\hat{S} - \check{D}|^{-\alpha}}$$

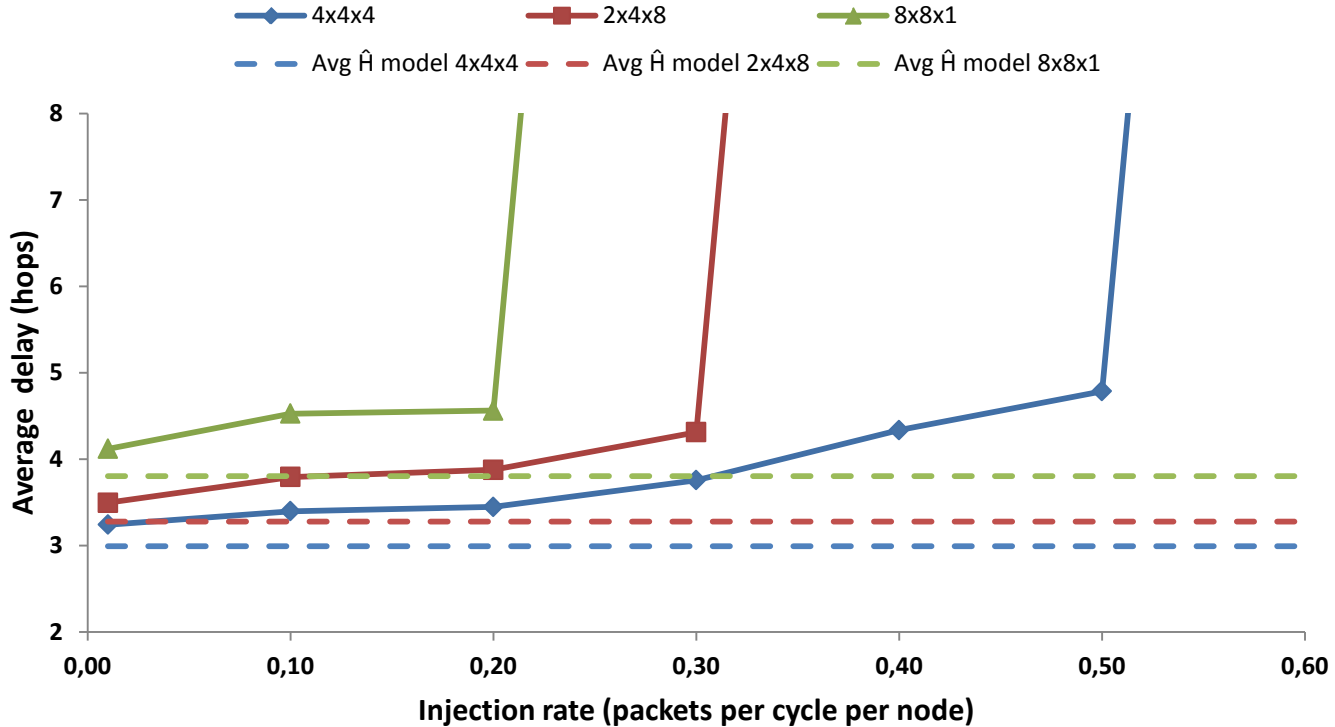
α -Model Localized Pattern

Effect of Localization Coefficient on Traffic Performance



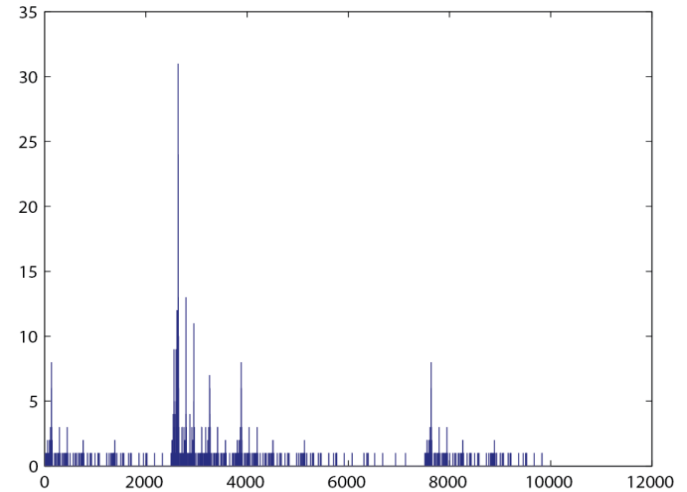
α -Model Localized Pattern

Avg. hop count with self-similar Local-alpha, bias, $\beta = 0.5$



Temporal Distribution (B-Model)

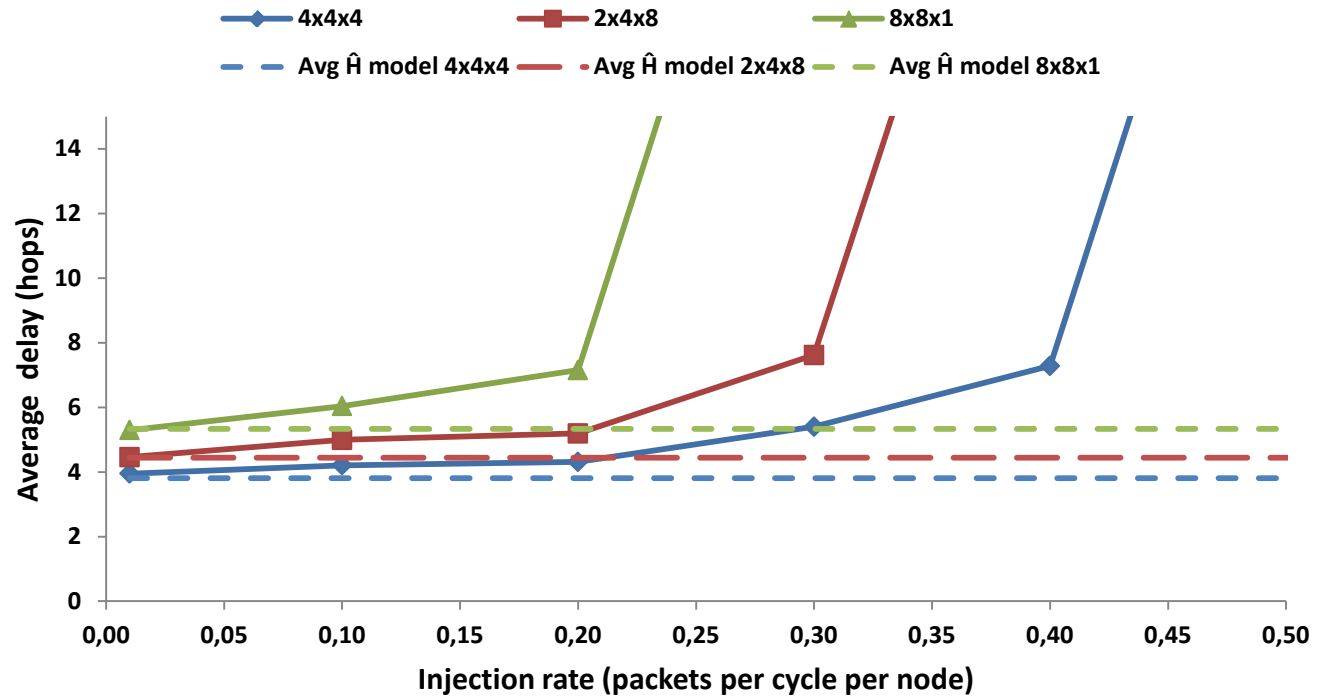
- Networks show a pattern of self-similarity in real traffic scenario.
- In a Bursty model, at any point, $x(i \cdot n/2^d)$, in a time series, the number of packets that a node injects to the network is expressed as a function of the bias, β , the division depth, d , and the injection rate, γ .



$$x\left(i \frac{n}{2^d}\right) = (\{\beta, 1 - \beta\})^d \left(\gamma n - \sum_{j=0}^{i-1} x\left(j \frac{n}{2^d}\right) \right)$$

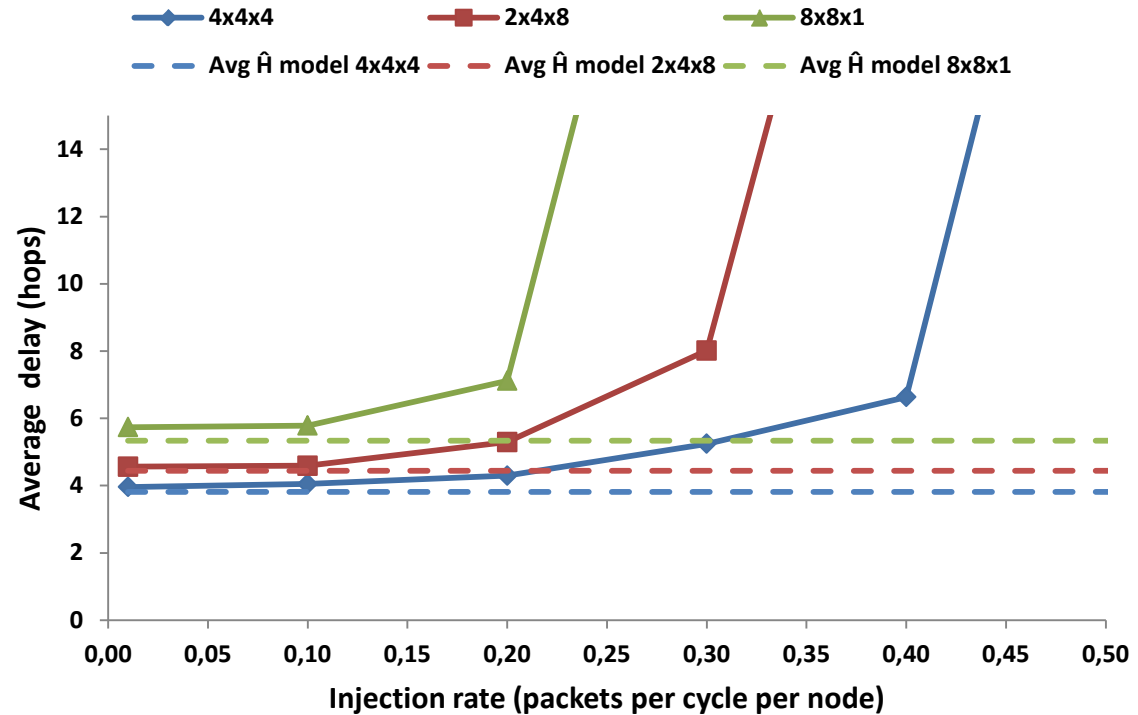
URT

Avg. hop count for URT, bias, $\beta = 0.5$



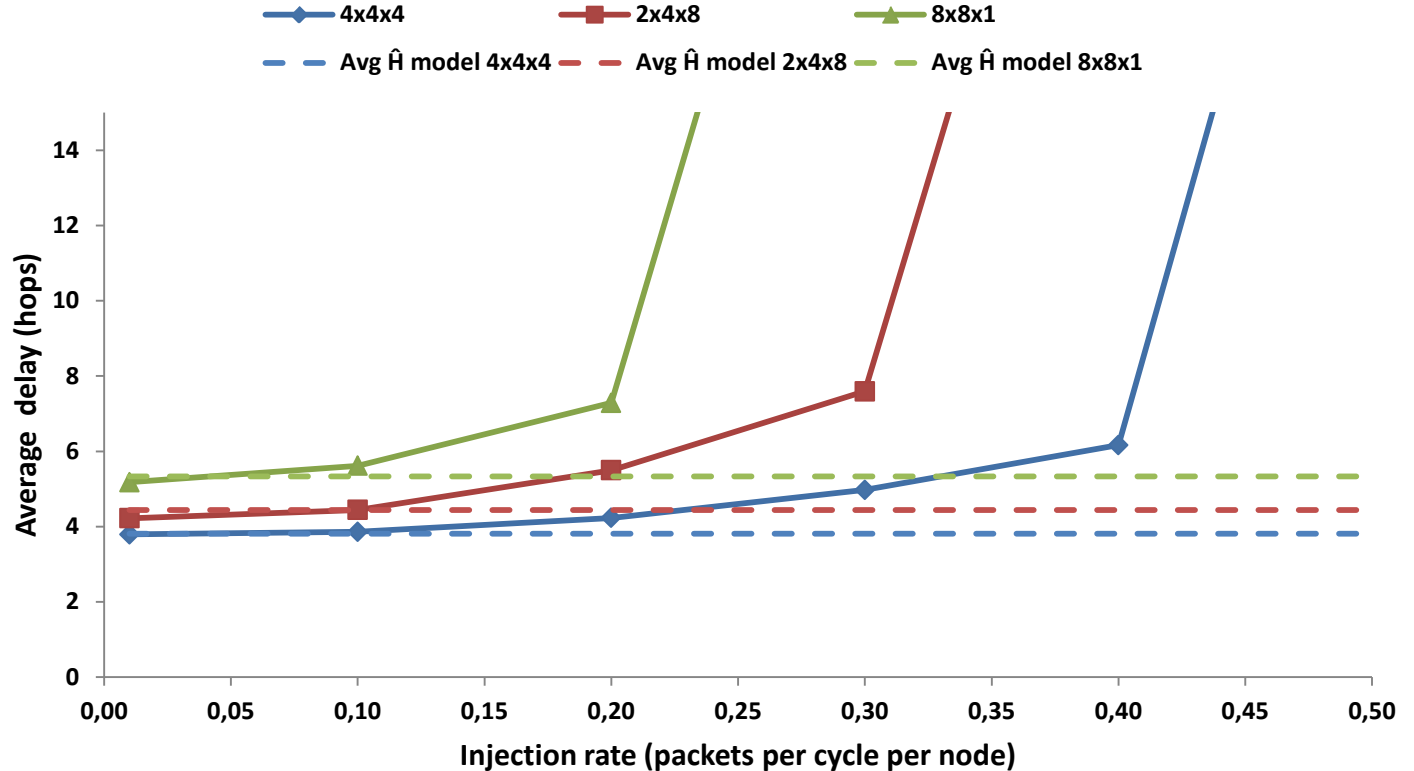
URT

Avg. hop count with self-similar URT, bias, $\beta = 0.3$



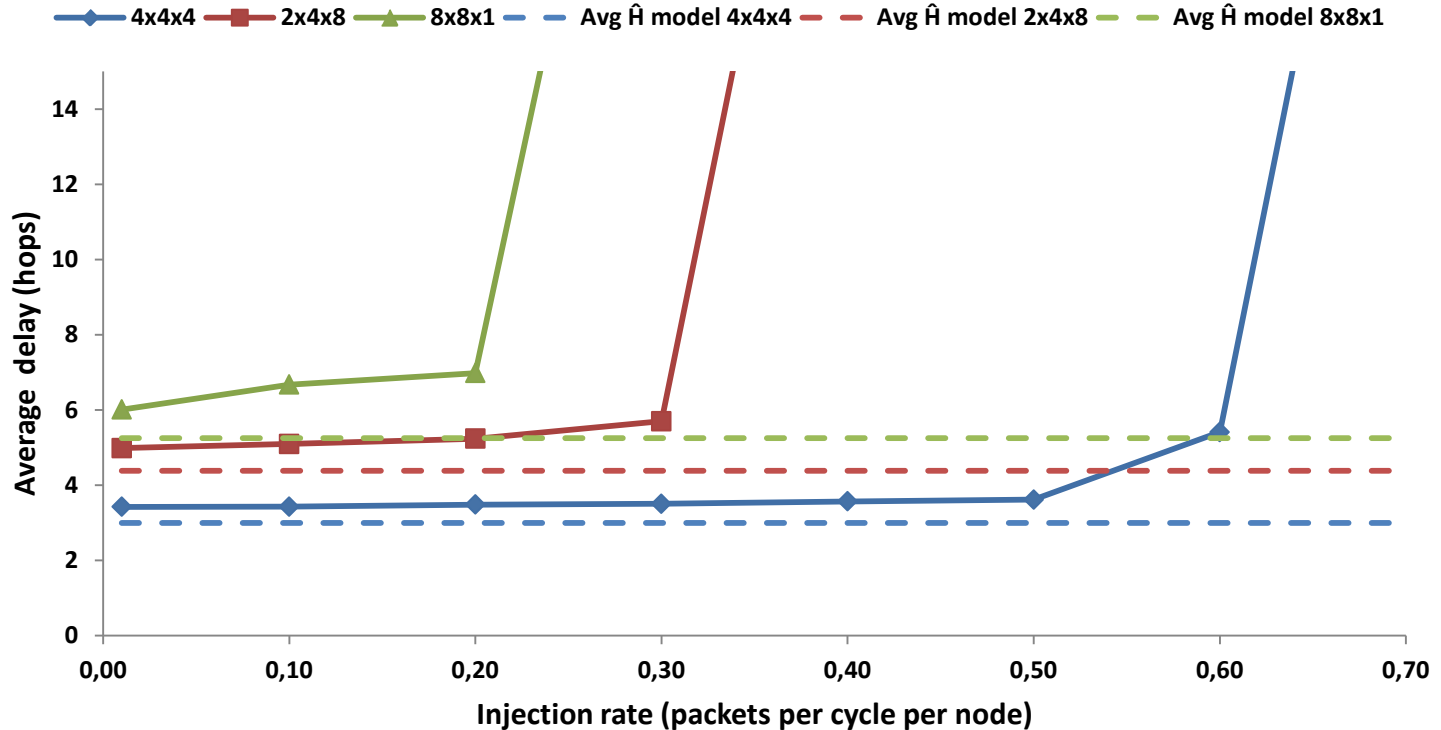
URT

Avg. hop count with self-similar URT, bias $\beta = 0.1$



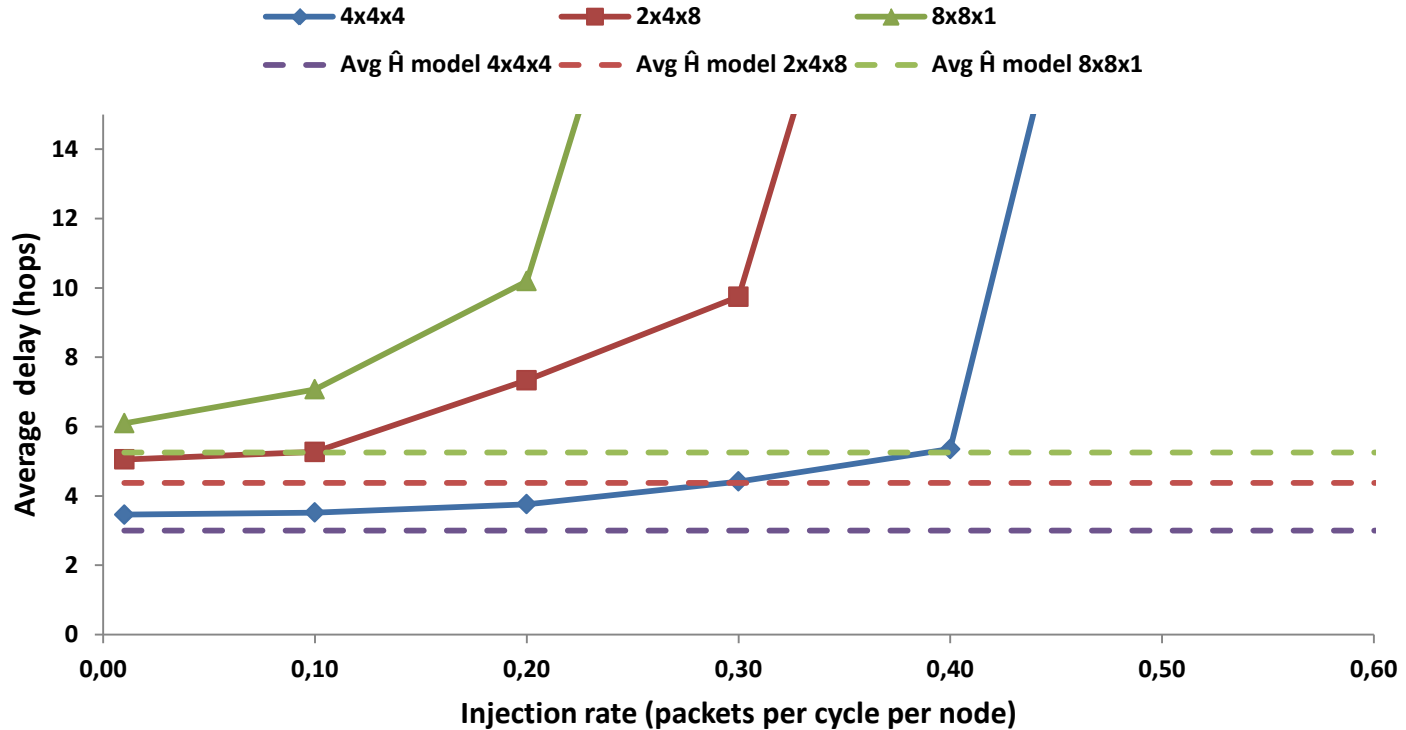
Bit-Reverse

Avg. hop count with Bit-reverse, bias $\beta = 0.5$



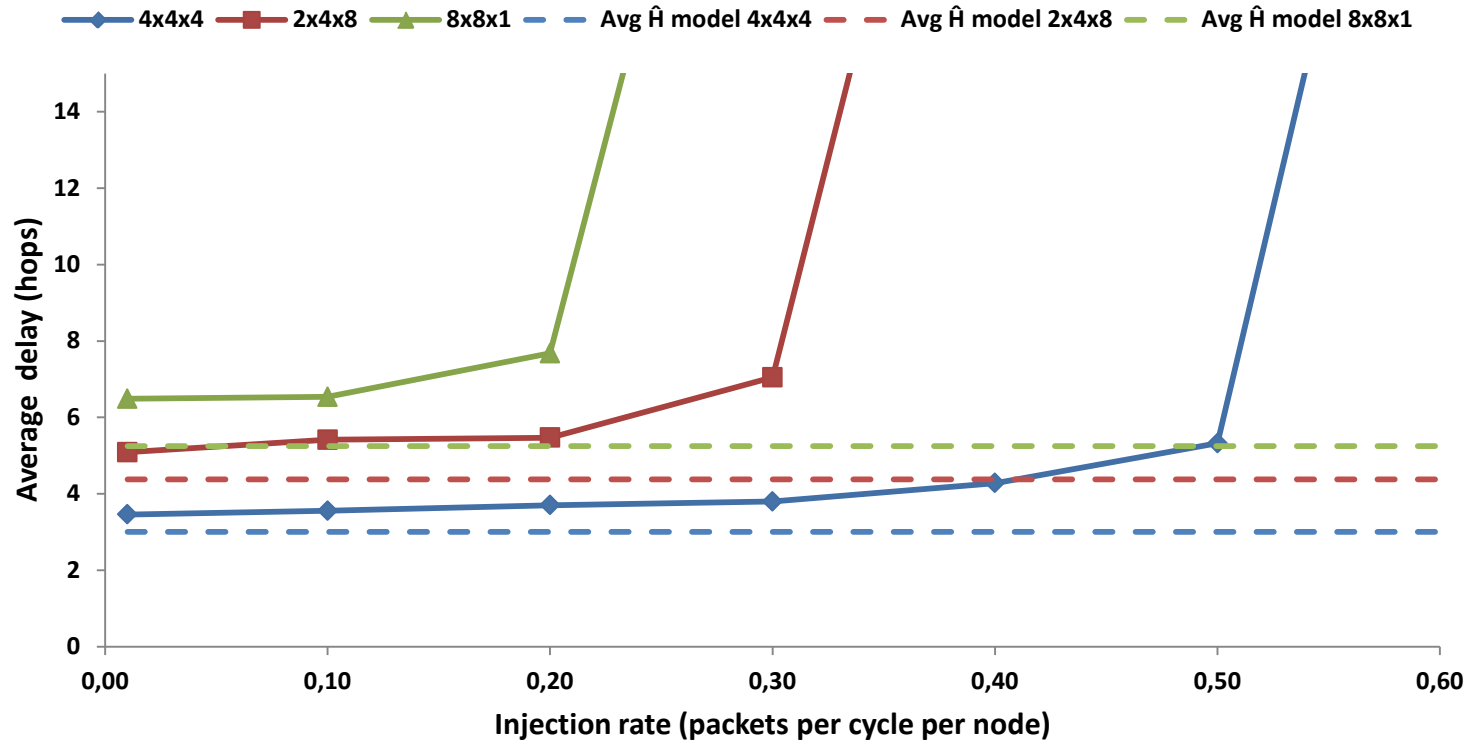
Bit-Reverse

Avg. hop count with self-similar Bit-reverse, bias, $\beta = 0.1$



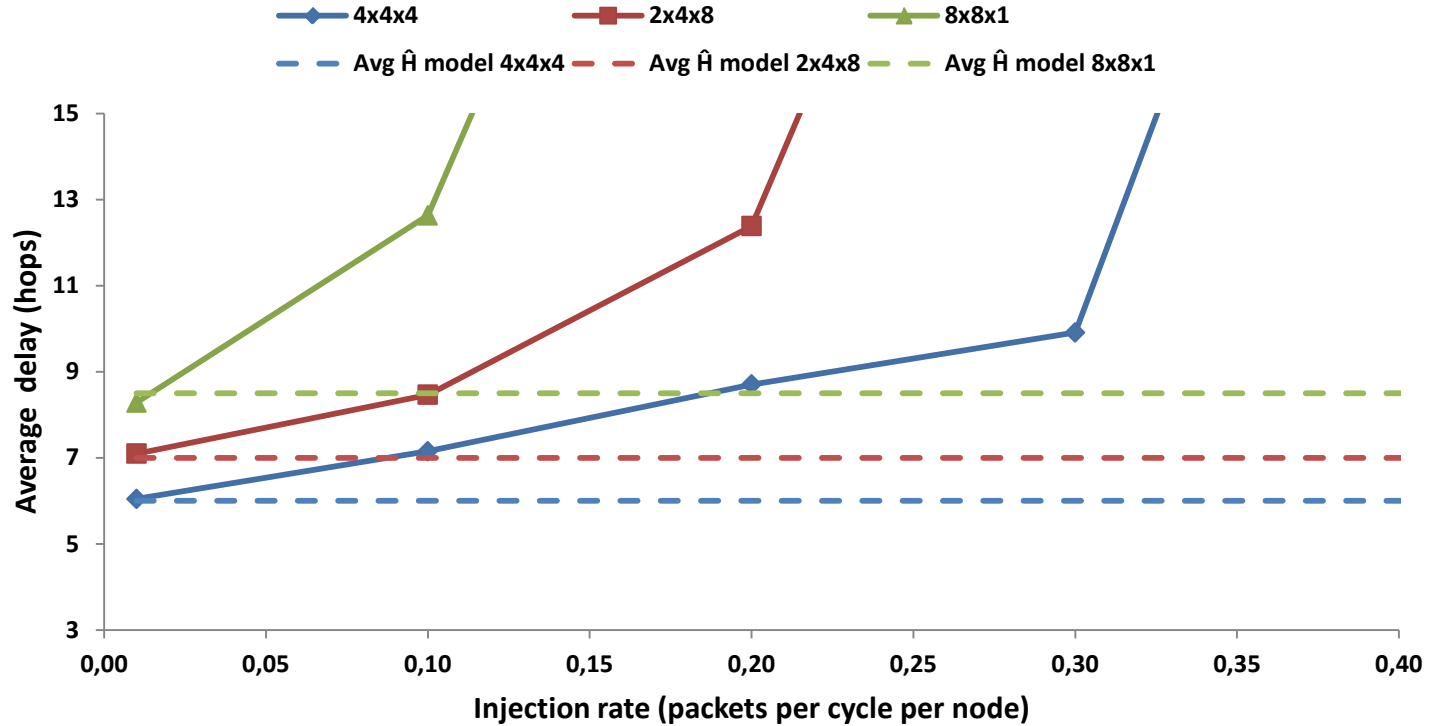
Bit-Reverse

Avg. hop count with self-similar Bit-reverse, bias, $\beta = 0.3$



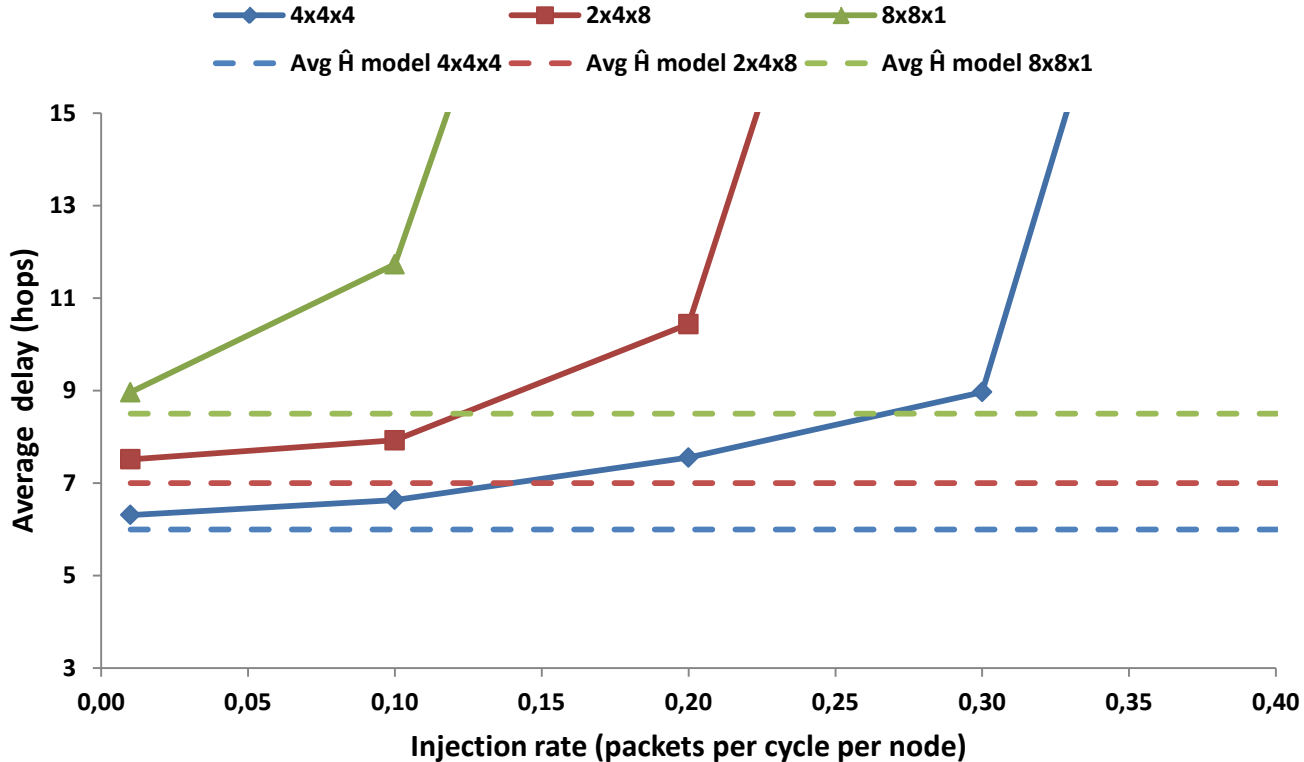
Bit-Complement

Avg. hop count with self-similar Bit-complement, bias, $\beta = 0.1$



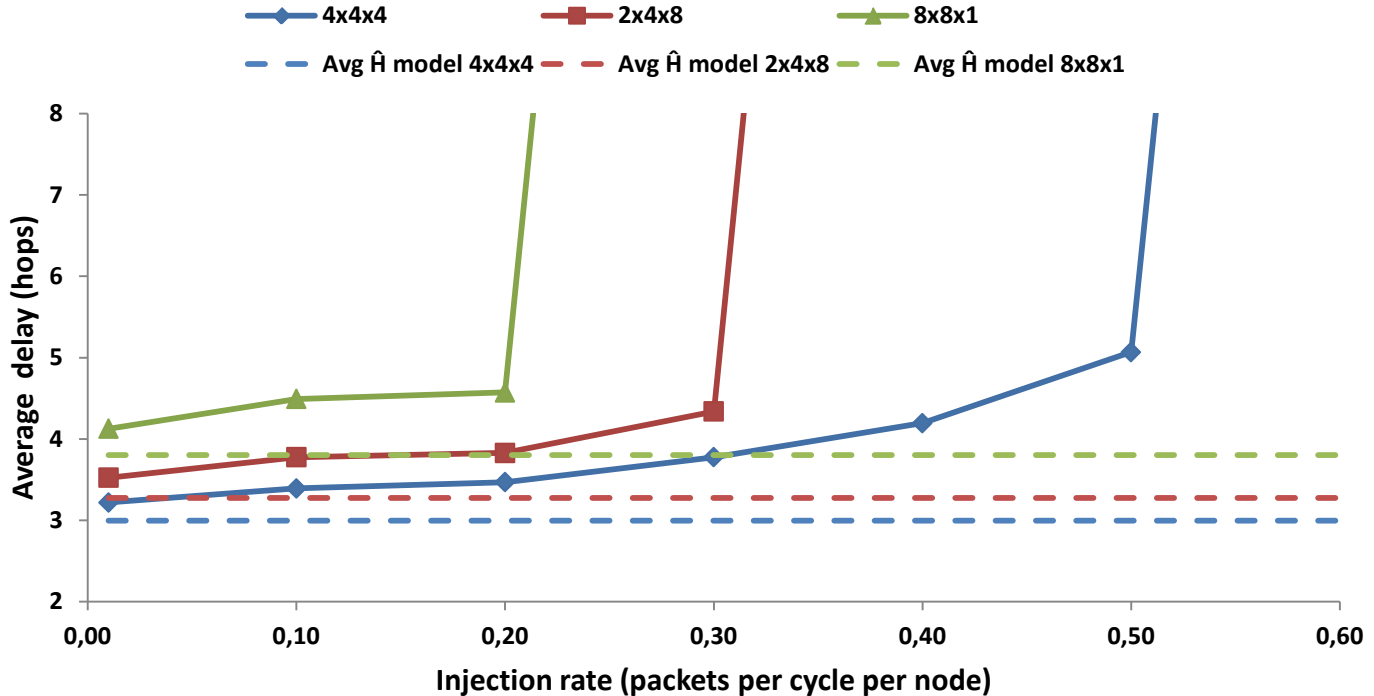
Bit-Complement

Avg. hop count with self-similar Bit-complement, bias, $\beta = 0.3$



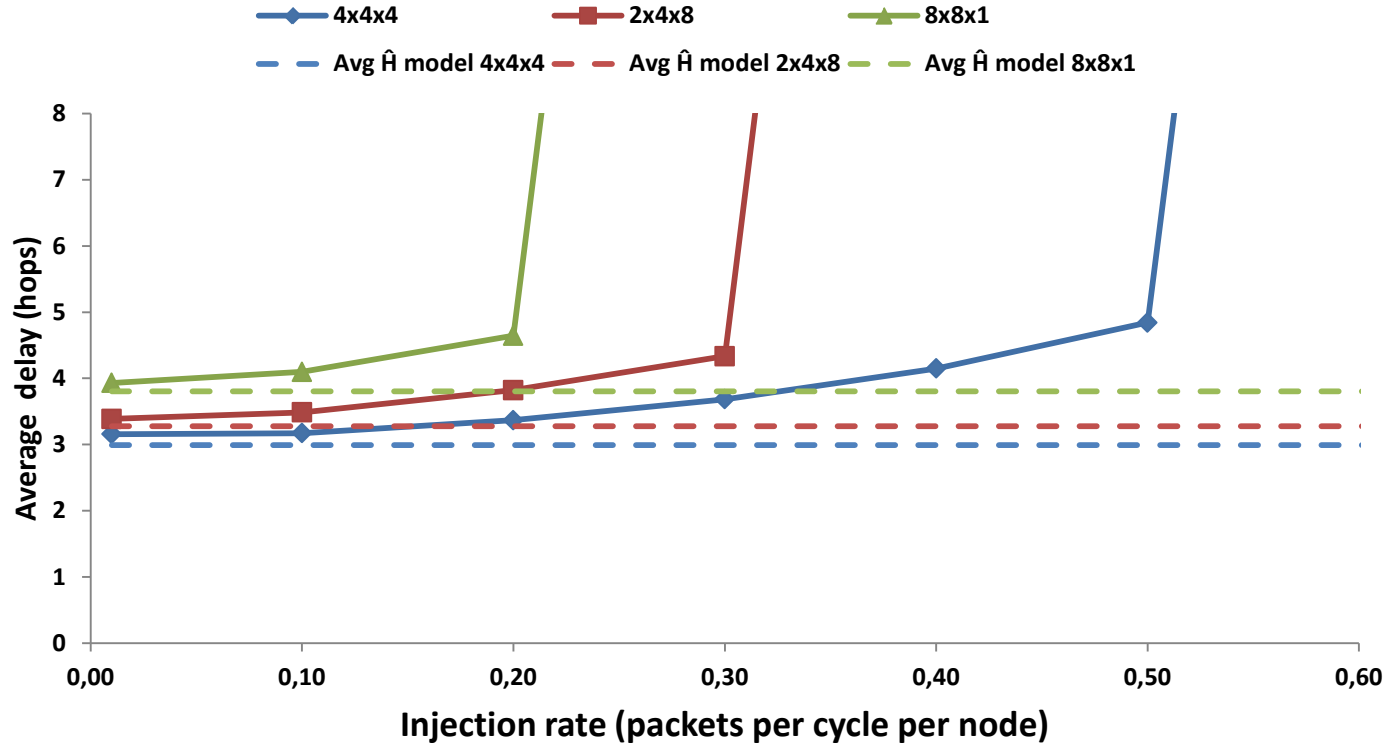
α -Model Localized Pattern

Avg. hop count with self-similar Local-alpha, bias, $\beta = 0.3$



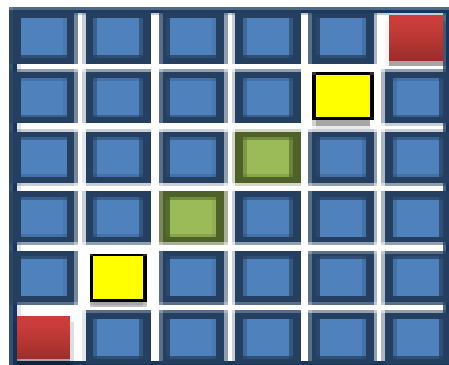
α -Model Localized Pattern

Avg. hop count with self-similar Local- α , bias, $\beta = 0.1$



Hotspots

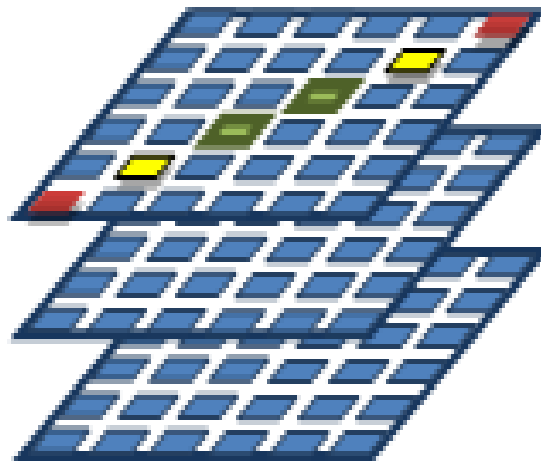
- The placement of hotspot nodes in the network affects the overall network performance.



HS1

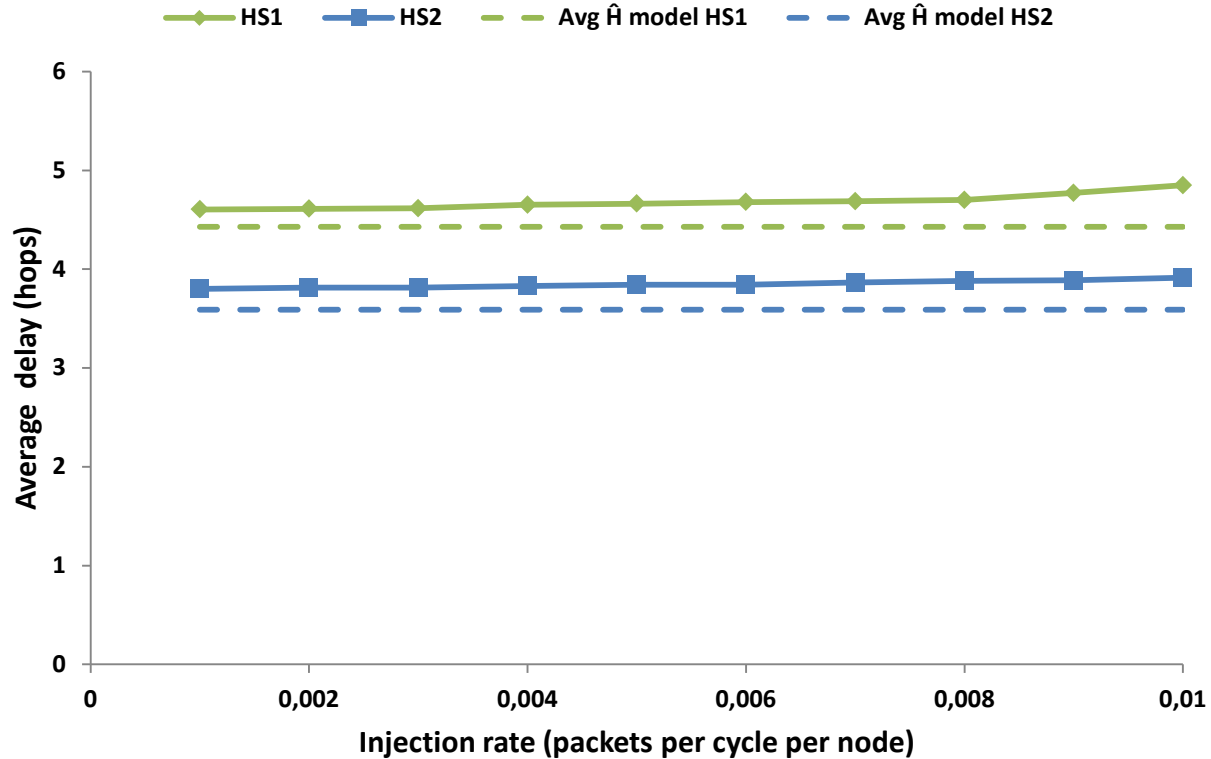
HS2

HS3



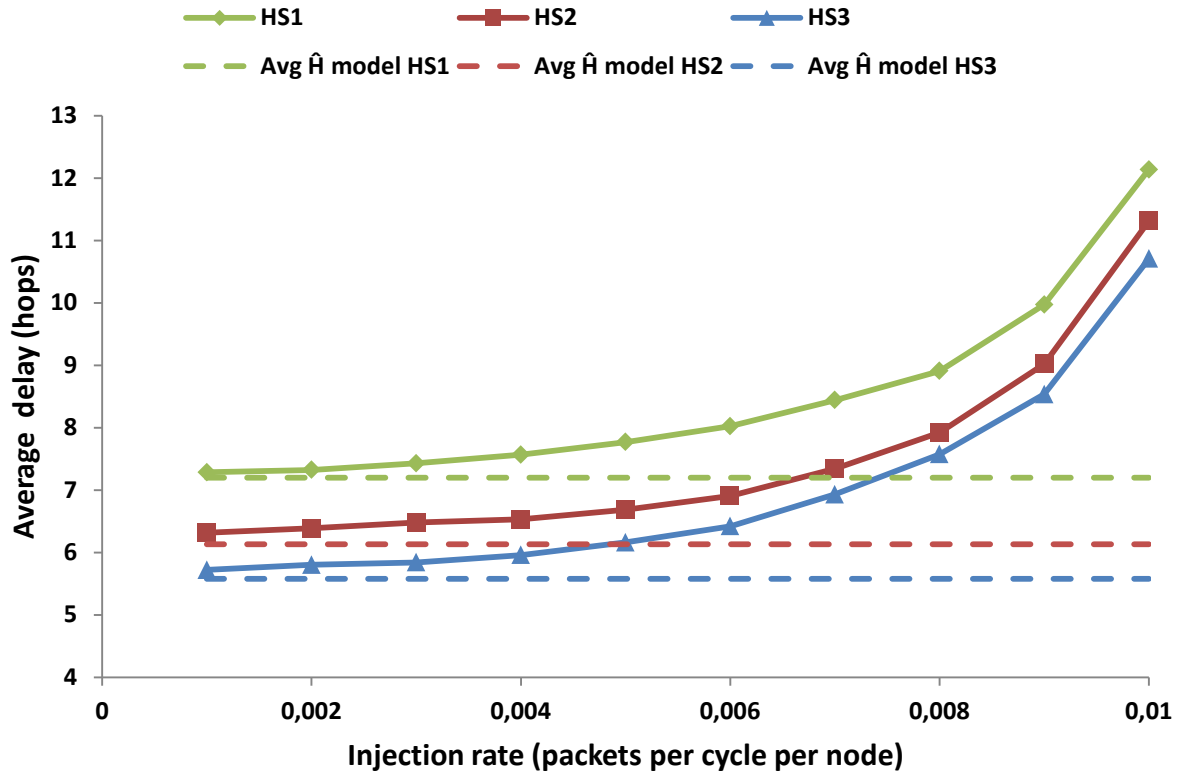
Hotspots

Avg. hop count, self-similar hotspot 4x4x4



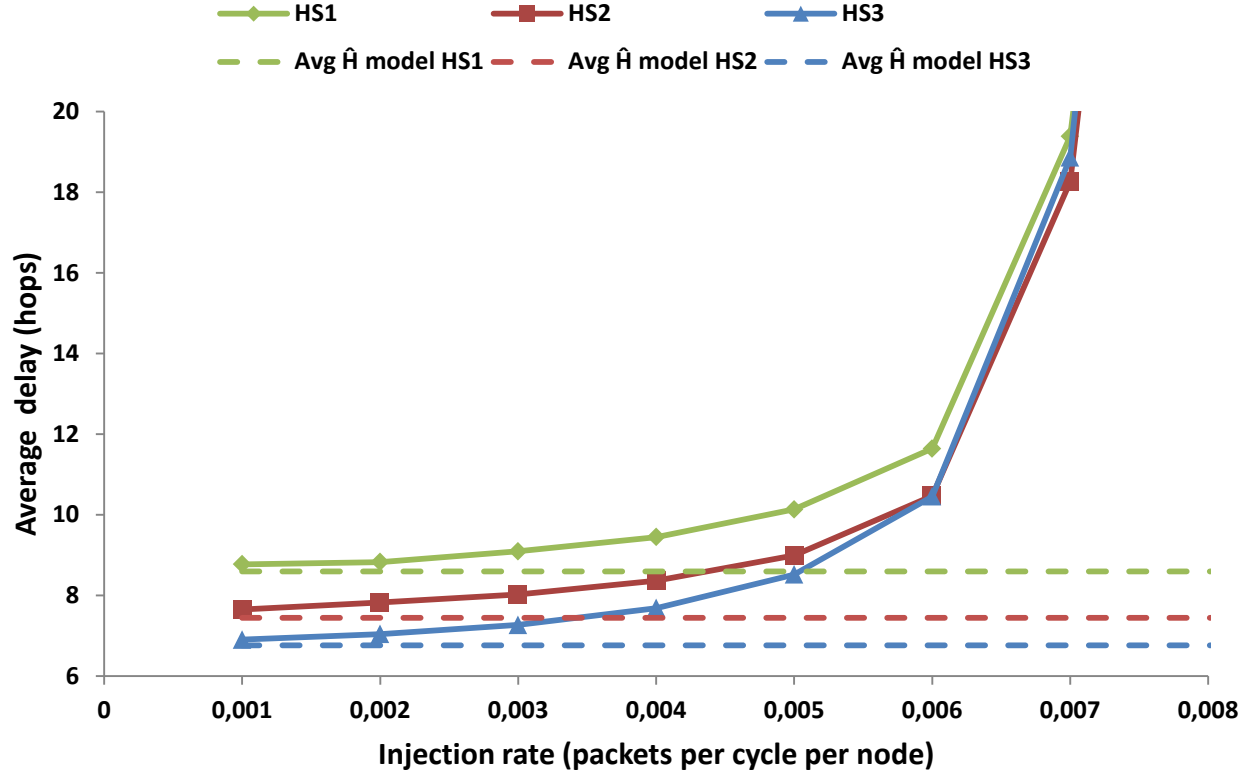
Hotspots

Avg. hop count, self-similar hotspot 6x6x6



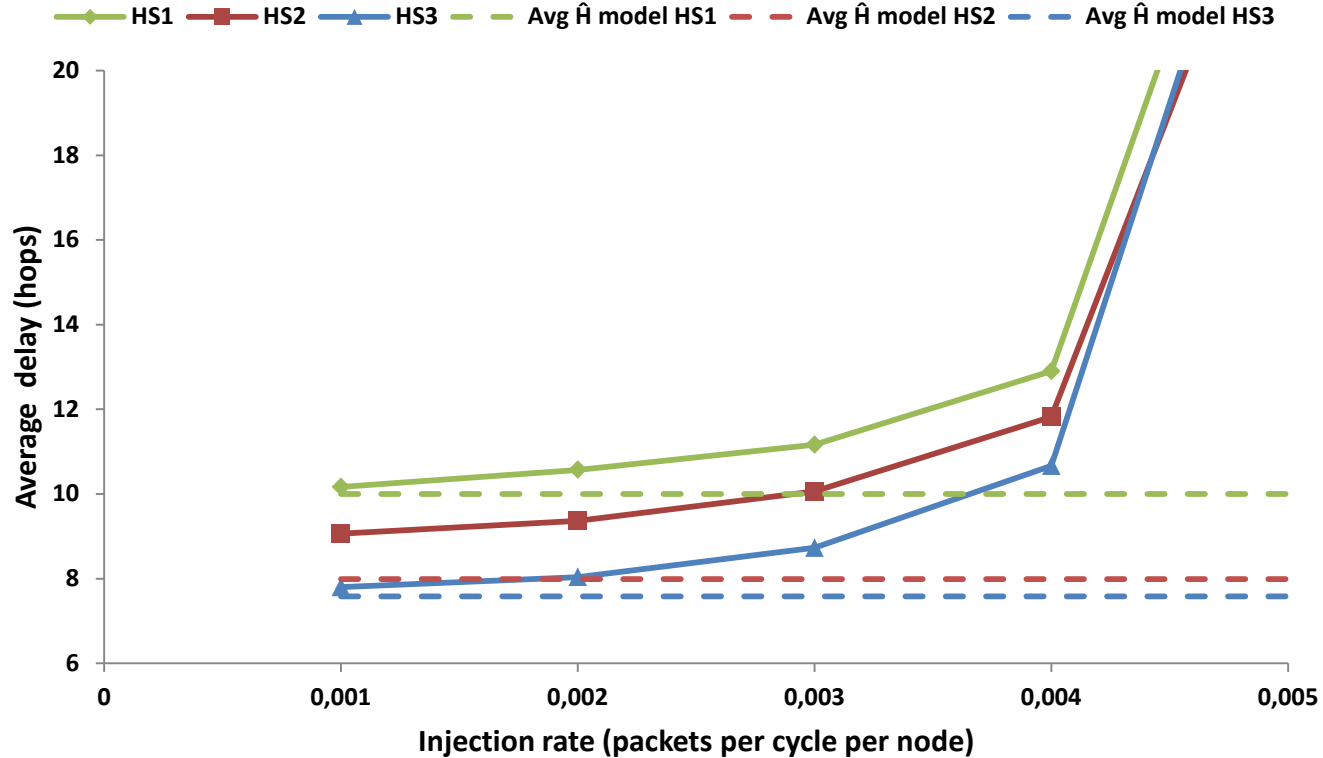
Hotspots

Avg. hop count, self-similar hotspot 7x7x7



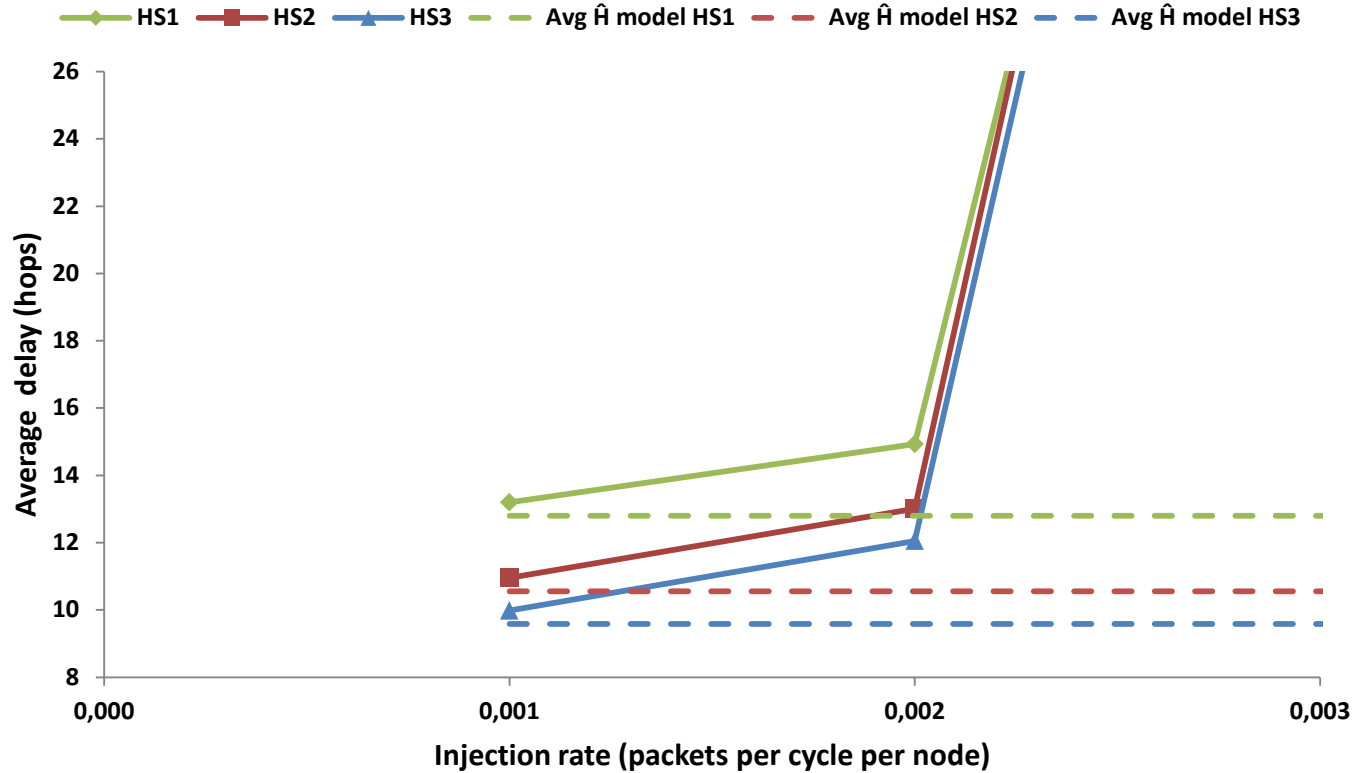
Hotspots

Avg. hop count, self-similar hotspot 8x8x8



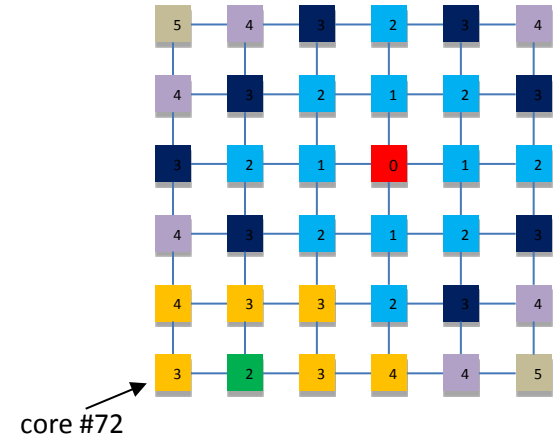
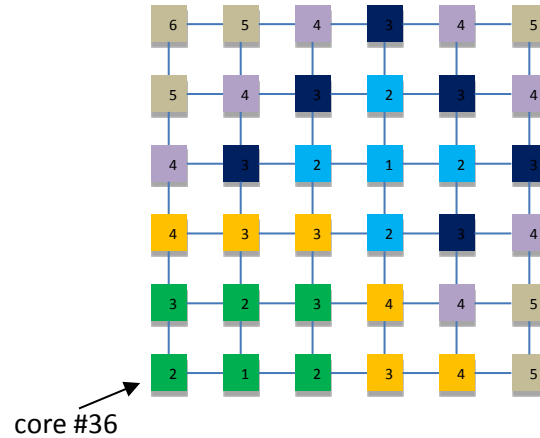
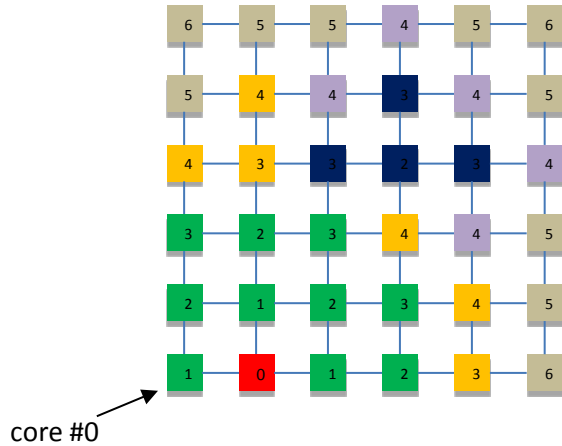
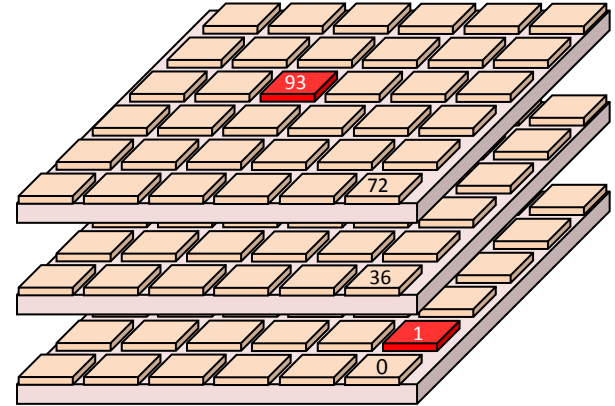
Hotspots

Avg. hop count, self-similar hotspot 10x10x10



Configuration

Number of cores	Memory access	Off-chip access	With other cores
A=18	7.14%	7.14%	85.71%
B=18	14.29%	14.29%	71.43%
C=18	21.43%	21.43%	57.14%
D=18	28.57%	28.57%	42.86%
E=18	35.71%	35.71%	28.57%
F=18	42.86%	42.86%	14.29%

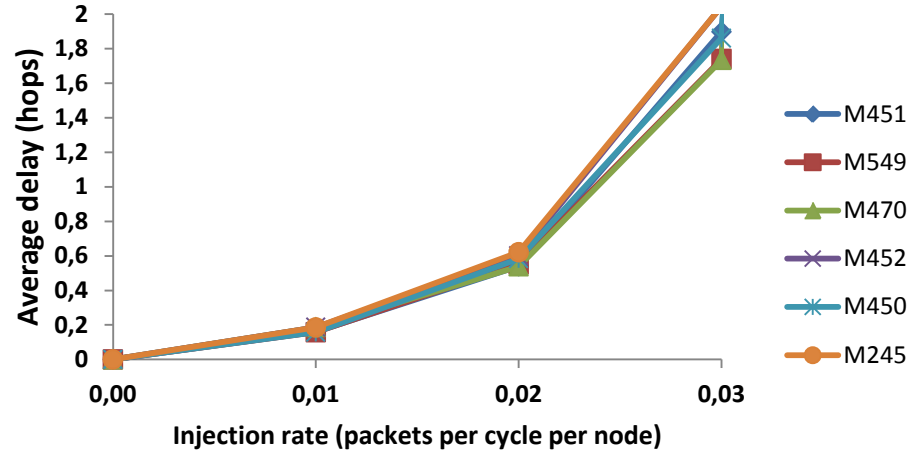


Configuration

➤ Output

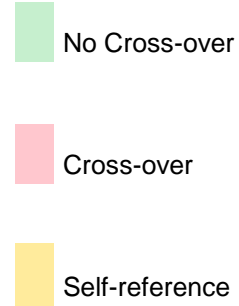
Injection-r	M451	M549	M470	M452	M450	M245
0,00	4,6041	4,6062	4,6085	4,6122	4,6451	4,888
0,01	4,82	4,85	4,83	4,84	4,90	5,13
0,02	5,18	5,24	5,18	5,29	5,29	5,56
0,03	6,53	6,42	6,39	6,71	6,58	6,96
0,04	21,09	25,25	28,83	27,51	28,24	29,01
0,05	118,52	132,25	131,28	132,17	135,10	162,46
0,06	306,91	307,70	325,68	319,67	330,98	334,17
0,07	448,99	448,91	449,40	450,26	456,25	455,13
0,08	537,41	541,05	543,97	543,24	548,50	535,00
0,09	603,57	587,39	607,01	588,38	602,23	603,87
0,10	627,37	604,70	623,91	611,32	605,12	644,24

Avg. hop count for M configurations



Cross-Overs

Injection-r	M451	M549	M470	M452	M450	M245
Model	4,6041	4,6062	4,6085	4,6122	4,6451	4,888
c(M451)	0,00%	-0,05%	-0,10%	-0,18%	-0,88%	-5,81%
c(M549)	0,05%	0,00%	-0,05%	-0,13%	-0,84%	-5,77%
c(M470)	0,10%	0,05%	0,00%	-0,08%	-0,79%	-5,72%
c(M452)	0,18%	0,13%	0,08%	0,00%	-0,71%	-5,64%
c(M450)	0,89%	0,84%	0,79%	0,71%	0,00%	-4,97%
c(M245)	6,17%	6,12%	6,06%	5,98%	5,23%	0,00%



$$\text{Difference \%} = \left| \frac{(z(M245) - z(M451))}{z(451)} \right| * 100$$

Configuration

- Cross-overs do happen when configurations are very similar.
- In all tested cases, there was never a cross over when the minimum difference in the zero-load prediction of two configuration $> \pm 0.13\%$.

Configuration to Minimize Average Delay

- The number of cores and network size is known.
- The preferred network dimensions in X,Y,Z is determined.
- The probabilities of each core sending packets to a hotspot node is estimated.
- The number and position of hotspot nodes is set.
- The cores are equally grouped into clusters based on their probabilities range.
- The Zero-load model finds configuration with the lowest average distance

Zero Load Predictive Model

- Static performance predictor
- Based on Network geometry and traffic pattern
- High fidelity observed in simulations
- Applications:
 - Fast and early design space exploration
 - Node placement
 - Task mapping
 - ...

Questions