

Effective Resource Model and Cost Scheme for Maze Routing in 3D Global Routing

Zirui Li¹, Jianwang Zhai^{1†}, Zixuan Li¹, Zhongdong Qi², Kang Zhao¹

¹Beijing University of Posts and Telecommunications ²Xidian University

{lzh_offical, zhaijw, gxydbc, zhaokang}@bupt.edu.cn, zdqi@xidian.edu.cn

Abstract—Routing is an essential step in the design closure of integrated circuits (IC) and has become the runtime bottleneck in the physical design flow of very large-scale integrated (VLSI) circuits. The maze routing approach, which is mostly used in the rip-up and reroute (RRR) stage of global routing, closely affects routing efficiency and quality. In this paper, we proposed an effective resource model and dynamic congestion sensitivity adjusting method for multi-level 3D maze routing, to improve the efficiency and quality of solutions for 3D global routing. The proposed method takes into account the resource distribution in the maze routing planning process, can optimize the composition of the solution space, and can flexibly adjust the congestion sensitivity. The experimental results show that after integrating the proposed model in a high-performance 3D global router, better routing quality can be achieved in shorter runtime.

Index Terms—3D Global Routing, Rip-up and Reroute, Resource Model, Multi-level 3D Maze Routing

I. INTRODUCTION

As the final process of physical design, routing is an essential step to VLSI design closure. In order to improve efficiency, modern routing flow is generally divided into global routing and detailed routing. Global routing is to find rough routes within the routing area under a series of constraints, aiming at the minimum cost while taking into account detailed routability. However, the initial routing might not meet all the requirements due to limited routing resources. In such cases, the global router needs to remove some existing connections and then re-find and re-establish better routes to meet the requirements, which is known as rip-up and reroute (RRR). As shown in Fig. 1, the runtime percentage of RRR in global routing is very large, and more efficient models and algorithms are urgently needed.

RRR is an iterative process adopts specialized algorithms and optimization techniques to minimize congestion, a factor which was inadequately considered in initial routing and probably cause physical violations. To obtain better routing results, maze routing [1] is one of the commonly operated methods in the RRR iterations. This method runs on a grid map representing the routing region and amount of routing resources. The goal is to find the shortest path under routing resource constraints.

Many competitive global routers [2]–[6] compress the 3D routing space [7] into 2D to handle VLSI circuits, which is

[†]Corresponding author. This work is supported in part by National Key R&D Program of China (2022YFB2901100) and Beijing Natural Science Foundation under Grant 4244107.

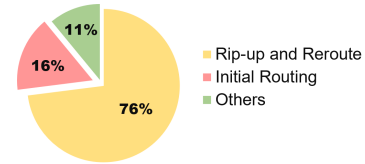


Fig. 1 Breakdown runtime results on ICCAD2019 benchmarks [11] by CUGR [9] given in [12].

effective in reducing routing runtime but inevitably sacrifices routing quality. By directly searching paths in 3D routing space, better routing quality can be achieved, and multi-level routing can be used to speed up the path search process. MGR [8] adopts a multi-level framework and two-stage 3D maze routing method to respectively plan and search for routes at coarse and fine granularity. It greatly balances the performance and runtime. CUGR [9] proposes a multi-level 3D maze routing algorithm with a probability-based cost scheme that better considers detailed routability, and applies it to the RRR iterations. Based on CUGR, FastGR [10] proposes a global routing framework accelerated for heterogeneous CPU-GPU systems to perform parallel routing from the architectural level to improve efficiency.

Existing 3D global routers optimize algorithms and frameworks from different perspectives to improve routing efficiency and quality, and have achieved remarkable results. However there is room for improvement in accurate modeling of routing resources and cost function in multi-level 3D maze routing. Taking account of distribution while evaluating the resources or more precise value setting of parameters may bring performance in runtime and solution quality.

In this work, we propose an effective resource model and congestion adjusting method, to pursue higher efficiency and better results. The main contributions of this paper include:

- A resource weighting method that takes into account the distribution of G-cells. Different weight models can be established for wire and via edges.
- A dynamic adjustment method to congestion sensitivity in cost scheme.
- Experiments show that using the multi-level 3D maze routing based on our model, the RRR speed and global routing efficiency are improved by 6.64% and 5.69%, respectively, and better routing quality is achieved.

The rest of the paper is organized as follows: Section II introduces the base model and algorithm process. Section III

TABLE I Notations

Symbol	Definition
u, v	G-cells
(u, v)	A wire edge or via edge with u and v as endpoints
$c(u, v)$	Capacity of edge (u, v)
$d(u, v)$	Demand of edge (u, v)
$r(u, v)$	Resource of edge (u, v)
$wl(u, v)$	Wire length of a wire edge (u, v)
$eo(u, v)$	Expectation of overflow of a wire edge (u, v)
$lg(u, v)$	Logistic coefficient of an edge (u, v)
uoc	Unit overflow cost
$slope$	An adjustable parameter of overflow sensitivity
s	The scale of coarsening

details the proposed methodology. Section IV gives the experimental results. Section V concludes the paper.

II. PRELIMINARIES

A. Problem Formulation

Global routing works on a set of horizontally and vertically evenly distributed G-cells. By considering each G-cell as a vertex $v \in V$ and connections between adjacent G-cells as an edge $e \in E$, the routing area can be formed as a grid graph $G(V, E)$. The connections between G-cells on the same layer and different layers are called wire edges and via edges, respectively. The capacity of edges is defined as the number of tracks going through. The number of tracks that remain and are occupied are called the edge's resource and demand. Finding a route connecting all the pins under design rules with minimum cost while meeting the capacity requirements is the main purpose of global routing.

Many modern global routers [5], [6], [9], [10], [13] divide their flow into two parts: the initial routing stage and the RRR iterations. The initial routes shape has been determined through pattern routing, and the nets with violations are dealt with in RRR iterations using maze routing. In addition, the modeling of the cost scheme also determines the quality of the solution.

B. Multi-level 3D Maze Routing

Multi-level 3D maze routing is an algorithm used for RRR. Compared to traditional maze routing methods, it introduces a technology called coarsening to narrow the 3D search space and divide RRR into two levels in different granularity.

Fig. 2 shows the procedure of multi-level 3D maze routing. First, blocks of multiple G-cells are compressed at the scale of s into coarsened cells which reconstruct the original grid graph into a coarsened grid graph. The resource of a coarsened cell is computed as the average resource of all the G-cells within it using Equation (1), in which i, j are the index of row and column in coarsened cell A.

$$R(A) = \frac{\sum_{i=0}^{s-1} \sum_{j=0}^{s-1} r(v_{ij})}{s^2}. \quad (1)$$

Then, coarse-grained maze route planning is performed on the coarsened grid graph to find the rough search area of the optimal route. In this way, the fine-grained maze router will find the minimum cost route within the planned area instead of directly searching the whole 3D graph.

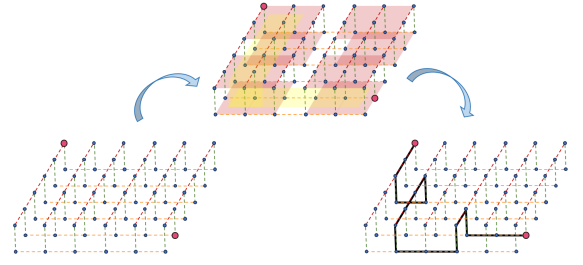


Fig. 2 The illustration of multi-level 3D maze routing.

C. Cost Scheme

Our proposed methodology is based on the competitive global router, CUGR [9], whose primary notations are listed in TABLE I. The cost scheme proposed by CUGR not only takes into account the congestion situation during the routing process, but also dynamically adjusts the sensitivity to congestion in different RRR iterations. Take the wire edge as an example, the cost function is:

$$\begin{aligned} cost_{w(u,v)} &= wl(u,v) + eo(u,v) \times lg(u,v), \\ eo(u,v) &= wl(u,v) \times \frac{d(u,v)}{c(u,v)} \times uoc, \\ lg(u,v) &= \frac{1}{1 + e^{(slope \times r(u,v))}}. \end{aligned} \quad (2)$$

The parameter $slope$ determines to what extent of resource scarcity lead to the rapid increase of the congestion cost. In other words, $slope$ determines the router's congestion sensitivity. Suppose i is the round of RRR iteration, CUGR update $slope$ in every RRR iteration by multiply it by 2^i to adjust the congestion sensitivity.

Defining the possibility of overflow as $[1 - r(u,v)/c(u,v)]$, the congestion weight $wcg_w(u,v)$ can be expressed as Equation (3). Reorganize the cost function as Equation (4) to more intuitively understand the composition and changing trend with congestion of the cost. The congestion cost weight is virtually 0 when the resources are sufficient. Whereas it will increase rapidly maximum to 0.5 on condition that resources are scarce. Thus how the parameter update method is designed to adjust congestion sensitivity is particularly important.

$$wcg_w(u,v) = uoc \times \left[1 - \frac{r(u,v)}{c(u,v)} \right] \times lg(u,v). \quad (3)$$

$$cost_w(u,v) = wl(u,v) \times [1 + wcg_w(u,v)]. \quad (4)$$

III. METHODOLOGY

A. Optimization of Multi-level 3D Maze Routing

The reduction of the solution space of the multi-level 3D maze routing mainly comes from the coarsening of the entire routing grid graph. After coarsening, maze route planning is performed to narrow the search space of the fine-grained maze routing process, which will improve efficiency at the cubic level. Therefore, the processing of the coarsened cells is of vital significance for accelerating the stage of RRR. As shown in Equation (1), CUGR directly compute the resource of a

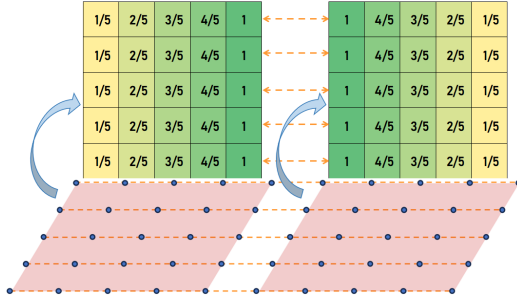


Fig. 3 Example of weight distribution in coarsened edges.

coarsened cell by the arithmetic mean value of resources of all the G-cells within. This method effectively quantifies the overall resource of a coarsened cell, yet it ignores the resource distribution of G-cells inside.

Since each layer can only be routed horizontally or vertically, the importance of resources varies among G-cells at different positions. Take the wire edge as an example for illustration. G-cells aggregate into coarsened cells and coarsened edges are formed by the interconnections between coarsened cells. The closer the G-cells within coarsened cells are to the connection, the more important its resources are. If the intersecting boundary resources of two adjacent coarsened cells are insufficient to form a continuous edge, regardless of the abundance of resources in their respective intermediate regions, they cannot establish a path. Therefore, the G-cell resource in each coarsened cell should be given corresponding weights according to its location.

A wire edge is split into several pairs of adjacent coarsening units, higher weights should be assigned to resources of G-cells near the connection. We use Equation (5) to calculate the weight $w(u_{ij})$ of resources in G-cell u_{ij} , assuming A and B is a pair of horizontally adjacent coarsened cells.

$$w(u_{ij}) = \begin{cases} \frac{x_{ij} - x_{low}}{s_x}, u_{ij} \in A \\ \frac{x_{high} - x_{ij}}{s_x}, u_{ij} \in B \end{cases}, \quad (5)$$

$$w(u_{ij}) = \begin{cases} \frac{y_{ij} - y_{low}}{s_y}, u_{ij} \in A \\ \frac{y_{high} - y_{ij}}{s_y}, u_{ij} \in B \end{cases}.$$

In layers routed horizontally, x_{high} and x_{low} represent the largest and smallest value of horizon abscissa in a pair of coarsened cells, respectively. The parameter s_x refers to the coarsening scale. Similarly, we use y_{high} , y_{low} , and s_y to represent the corresponding ordinate in the layer routed vertically. An example resource weight distribution of a pair of adjacent coarsened cells is shown in Fig. 3.

The weight of each G-cell's resource changed from 1 to $w(u_{ij})$ after calculating the weighted sum of all G-cells' resources along the wire edge. Therefore, for s G-cells, the sum of weight changes from $s \times 1$ to $\frac{1}{s} + \frac{2}{s} + \dots + \frac{s-1}{s} + 1 = \frac{s+1}{2}$. It is necessary to multiply by $\frac{2s}{s+1}$ after the weighted sum, in order to ensure that the original relative weight of the wire edge and the via edge remains unchanged. The resources of coarsened cells in wire edges is computed as:

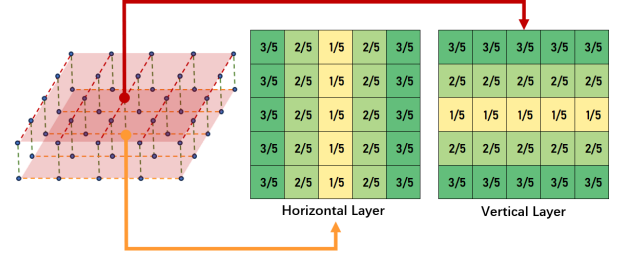


Fig. 4 Example of weight distribution in coarsened vias.

$$R(A) = \sum_{i=0}^{s-1} \sum_{j=0}^{s-1} r(u_{ij}) \times w(u_{ij}) \times \frac{2s}{s+1}. \quad (6)$$

For vias, higher weight should be assigned to the resources near the edge area. Suppose that A and B are a pair of coarsened cells in a via, respectively located in layers routed horizontally and vertically. Using (x_{mid}, y_{mid}) to represent the center coordinates of a coarsened cell, the weight $w(u_{ij})$ of resources in G-cell u_{ij} is calculated using the formula in Equation (7).

$$w(u_{ij}) = \begin{cases} \frac{|x_{ij} - x_{mid}|}{s_x}, u_{ij} \in A \\ \frac{|y_{ij} - y_{mid}|}{s_y}, u_{ij} \in B \end{cases}. \quad (7)$$

Fig. 4 shows an example resource weight distribution of a pair of coarsened cells in a via. Similarly, the weight sum of s G-cells in a coarsened via changes from $s \times 1$ to $\frac{s+1}{2} + \dots + \frac{2}{s} + \frac{1}{s} + \frac{2}{s} + \dots + \frac{s-1}{s} = \frac{s+1}{2}$. Equation (8) is used to compute the resources of coarsened cells in via edges.

$$R(A) = \sum_{i=0}^{s-1} \sum_{j=0}^{s-1} r(u_{ij}) \times w(u_{ij}) \times \frac{2s}{s+1}. \quad (8)$$

B. Dynamic Adjustment Method to Congestion Sensitivity

We propose a more accurate parameter update method, which is used to adjust the sensitivity to congestion of different RRR iterations, so as to obtain a better solution. To ensure that the congestion cost starts increasingly weighted from the same level with the reduction of resources in each iteration of RRR, *slope* and resource $r(u, v)$ need to satisfy:

$$slope \times r(u, v) = \lambda. \quad (9)$$

We define λ as a given constant that determines the level at which the congestion cost weights begin to increase. The resource is actually the number of unused tracks, therefore we select integer resource to calculate the value of *slope*. The parameter i_l is the number of RRR limits, and i is the number of RRR iterations. For the i^{th} iteration, the weight of congestion cost $wcg_w(u, v)$ should start increase when the resource $r(u, v) = i_l - i$.

$$slope = \lambda / (i_l - i). \quad (10)$$

Equation (10) is used to calculate *slope*, and the starting points in each iteration that $wcg_w(u, v)$ increasing are unified.

TABLE II Comparison of Decomposed GR Scores and Runtime

Design	Wirelength		# Vias		# Shorts		Score		RRR Runtime (s)		Overall Runtime (s)	
	CUGR	Ours	CUGR	Ours	CUGR	Ours	CUGR	Ours	CUGR	Ours	CUGR	Ours
18test5	27,007,300	27,009,300	855,453	854,563	0	0	16,925,500	16,922,900	21.143	18.660	34.943	31.457
18test5m	27,871,000	27,856,900	803,626	798,200	3,141	3,204	18,720,500	18,723,300	38.966	35.270	44.368	40.922
18test8	64,347,700	64,352,800	2,175,330	2,173,670	0	0	40,875,200	40,871,000	127.779	122.347	151.326	146.208
18test8m	64,700,300	64,574,100	1,955,550	1,950,320	5,527	5,738	42,935,600	42,957,100	151.379	139.463	168.172	156.627
18test10	66,772,100	66,748,300	2,309,520	2,306,100	0	0	42,624,100	42,598,500	143.817	139.102	167.954	163.016
18test10m	70,863,900	70,752,300	20,780,700	2,070,720	652	566	44,070,200	43,942,000	206.627	194.642	223.315	210.950
19test7	118,717,000	118,724,000	3,124,710	3,120,790	0	0	71,857,300	71,845,100	344.578	325.069	392.081	373.959
19test7m	106,995,000	107,054,000	30,713,500	30,560,300	4,367	4,284	67,966,200	67,893,100	209.886	203.209	249.219	236.754
19test8	181,899,000	181,883,000	57,481,700	57,452,900	0	0	113,942,000	113,922,000	178.908	160.141	249.465	230.910
19test8m	177,693,000	177,569,000	55,993,700	55,786,900	6,479	6,176	114,484,000	114,187,000	444.648	417.303	503.840	477.327
19tes9	181,899,000	181,883,000	57,481,700	57,452,900	0	0	175,532,000	175,496,000	241.538	220.801	365.717	343.861
19test9m	267,648,000	267,471,000	9,342,080	9,296,560	3,741	3,749	173,062,000	172,796,000	408.035	395.224	533.633	514.578
Average	-0.04%		-0.27%		-2.27%		-0.07%		-6.64%		-5.69%	

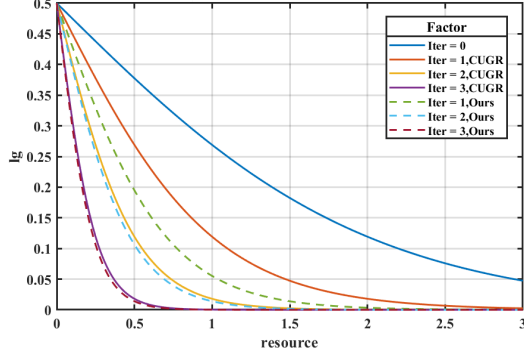


Fig. 5 Logistic curve in different iterations.

When $i_l = 4$, the logistic function $lg(u, v)$'s curves in each iteration are shown in Fig. 5. Compared with CUGR, our first RRR iteration reduces the sensitivity to a certain extent. Therefore, lower-cost path will be chosen in the first RRR iteration. It also can be seen that the sensitivity to congestion in the last two RRR iterations is almost unchanged. Consequently, the congestion degree of the routing result will not change much compared to the original version.

IV. EVALUATION

A. Experimental Setting

Our methodology is implemented in C++ based on the open-source global router CUGR [9]. Experiments are conducted on a Linux platform with Intel i9-13905H @5.4GHz CPU and 32GB memory. To evaluate the effectiveness of our methods, we adopted ICCAD2019 benchmarks [11] to estimate the performance and compare the routing efficiency and quality with CUGR.

A weighted sum score s is used to evaluate the quality of the global routing results, which is defined as Equation (11), where W means the wirelength, V is the number of vias, and S is the number of short violations, α , β , and γ are weighting factors set to 0.5, 4, and 500, according to [11].

$$s = \alpha W + \beta V + \gamma S. \quad (11)$$

B. Efficiency of RRR and Global Routing

As an important stage of global routing, the execution efficiency of RRR iterations directly affects the routing results.

Since our approach focuses on multi-layer 3D maze routing, the efficiency of RRR iterations is compared first. To measure the improvement in efficiency more fairly and accurately, each benchmark is run five times and averaged in the same environment. As shown in TABLE II, after our optimization the router is 6.64% faster in the RRR stage than CUGR, and the overall runtime of global routing is reduced by 5.69%. This is mainly due to our improvement to multi-level 3D maze routing. Since the spatial distribution of resources is fully considered in the RRR stage, the elimination of suboptimal options performs faster. Consequently, the solution space is optimized to a certain extent which leads to less runtime.

C. Quality of Global Routing

TABLE II also gives the comparison for the global routing results of all benchmarks. Overall, the wirelength is shortened by 0.04%, and the number of vias is reduced by 0.27%, which is mainly due to our optimization of the cost scheme. The maze router is given more routing options with lower costs because the weight of congestion cost in the first RRR iteration is decreased. Therefore, the wirelength and the number of vias are reduced to a certain extent. In the remaining iterations, the sensitivity of congestion is close to CUGR, which ensures minimal short violations due to congestion. As a result, the overall short violations are reduced by 2.29% though in several benchmarks the number of shorts increases. Finally, the average score is reduced by 0.07%.

It is worth mentioning that although the experiments are conducted based on CUGR, our proposed method is also effective for other global router based on multi-level 3D maze routing, such as FastGR [10].

V. CONCLUSION

In this work, we propose an effective resource model for the RRR iterations of global routing. First, the multi-level 3D maze routing is optimized by weighing resources of G-cells according to their distributions. Then, we redesign the parameter's updating method in the cost scheme to better adjust the sensitivity to congestion. As a result, our method achieves both routing efficiency and quality improvements on competitive global routers.

REFERENCES

- [1] J. Soukup, "Fast Maze Router," in *ACM/IEEE Design Automation Conference (DAC)*. IEEE Computer Society, 1978, pp. 100–101.
- [2] Y. Xu, Y. Zhang, and C. Chu, "FastRoute 4.0: Global router with efficient via minimization," in *IEEE/ACM Asia and South Pacific Design Automation Conference (ASPDAC)*, 2009, pp. 576–581.
- [3] W.-H. Liu, W.-C. Kao, Y.-L. Li, and K.-Y. Chao, "NCTU-GR 2.0: Multithreaded Collision-Aware Global Routing With Bounded-Length Maze Routing," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, vol. 32, no. 5, pp. 709–722, 2013.
- [4] H.-Y. Chen, C.-H. Hsu, and Y.-W. Chang, "High-performance global routing with fast overflow reduction," in *IEEE/ACM Asia and South Pacific Design Automation Conference (ASPDAC)*, 2009, pp. 582–587.
- [5] M. M. Ozdal and M. D. F. Wong, "Archer: A History-Based Global Routing Algorithm," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, vol. 28, no. 4, pp. 528–540, 2009.
- [6] Y.-J. Chang, Y.-T. Lee, and T.-C. Wang, "NTHU-Route 2.0: A fast and stable global router," in *IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, 2008, pp. 338–343.
- [7] J. A. Roy and I. L. Markov, "High-Performance Routing at the Nanometer Scale," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, vol. 27, no. 6, pp. 1066–1077, 2008.
- [8] Y. Xu and C. Chu, "MGR: Multi-level Global Router," in *IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, 2011, pp. 250–255.
- [9] J. Liu, C.-W. Pui, F. Wang, and E. F. Y. Young, "CUGR: Detailed-Routability-Driven 3D Global Routing with Probabilistic Resource Modelx," in *ACM/IEEE Design Automation Conference (DAC)*, 2020, pp. 1–6.
- [10] S. Liu, Y. Pu, P. Liao, H. Wu, R. Zhang, Z. Chen, W. Lv, Y. Lin, and B. Yu, "FastGR: Global Routing on CPU–GPU With Heterogeneous Task Graph Scheduler," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, vol. 42, no. 7, pp. 2317–2330, 2023.
- [11] D. Sergei, V. Alexander, W. Lutong, and X. Bangqi, "2019 CAD Contest: LEF/DEF Based Global Routing," Jan 2019.
- [12] J. Liu and E. F. Young, "EDGE: Efficient DAG-based Global Routing Engine," in *ACM/IEEE Design Automation Conference (DAC)*, 2023, pp. 1–6.
- [13] Y. Han, D. M. Ancajas, K. Chakraborty, and S. Roy, "Exploring High-Throughput Computing Paradigm for Global Routing," *IEEE Transactions on Very Large Scale Integration Systems (TVLSI)*, vol. 22, no. 1, pp. 155–167, 2014.