Incremental Layer Assignment for Critical Path Timing

Derong Liu, Bei Yu † , Salim Chowdhury ‡ and David Z. Pan

ECE Dept., University of Texas at Austin; [†]CSE Dept., Chinese University of Hong Kong; [‡]Oracle Corp.

Motivation and Problem

► Layer assignment assigns segments to metal layers after 2-D global routing.

Figure: General Global Routing Flow.

Wires and vias on top layers are wider and much less resistive.

Figure: The cross section of IC interconnection stack in advanced technology nodes. The normalized resistance values of different metal layers are listed in the table (source: [1] Hsu. et al. ICCAD'2014 [2] Yu. et al. ICCAD'2015).

▶ Segment delay depends on its assigned layer and its downstream capacitance, while via delay depends on the connecting two segments and their assigned layers.

Timing critical segments prefer higher metal layers.

Non-Critical Nets: n1 n2 ; Critical Net: n3

Figure: Incremental layer assignment example, (a) initial assignment, (b) assignment after timing optimization (source: [2] Yu. et al. ICCAD'2015).

Problem: Critical Path Layer Assignment

 $\mathbf{X}_{ii} = 1$; *j* $\forall i \in S(Nc) \quad j \in L$ Edge capacity constraint:

Given a 3-D grid graph, edge and layer information, initial routing and layer assignment, and set of critical nets, layer assignment re-assigns layers among critical and non-critical nets onto layers in order to minimize their maximum path timing and satisfy the edge capacity constraints.

Model Description

Elmore Delay Model

Figure: Net critical path timing considering segment and via delays.

If Uniform division by $K \times K$ may lead to unbalanced resource allocation. ► Self-adaptive partition provides similar number of segments in each partition. \triangleright Through multi-threading, each thread deals with a workload in a well-balanced manner.

Capacity Constraints

Figure: Example illustrations of capacity constraints, (a) edge capacity, (b) via capacity.

 \triangleright Edge capacity constraint – the maximum number of routing tracks along the edge. \triangleright Via capacity constraint – the maximum number of allowable vias in this grid.

CPLA Algorithms

ILP Formulation

Binary variables: *xij* represents segment *i* assigned on layer *j*; *yijpq* represents the via connecting segment *i* and *p* from layer *j* to *q*, which is equal to the product of *xij* and *xpq*.

Constraints Each released segment has to be assigned.

$$
\sum_{i\in\mathcal{C}(a)} x_{ij} \leq cap_e(j); \qquad \forall e\in E
$$

i∈*S*(*e*) Via capacity constraint:

> \sum (*i*,*p*)∈*Sx*(*Nc*) $y_{ijpq} + n_v \cdot (x_{ij} + x_{pq}) \leq cap_g(l), \ \forall l,j < l < q.$

Self-adaptive Quadruple Partition

K × *K* **Partition supports parallel scheme.**

- Incremental layer assignment can be solved in each partition separately.
- \triangleright Multithreading is applied to provide further speed-up.
- \triangleright The most recent updated results by peer threads can be taken into accounts.

Figure: Example of grid partition, (a) nets partition, (b) routing density for benchmark adaptec1 by NCTU-GR.

> **Figu** Figure: Performance comparison with ILP, (a) *Avg*(*Tcp*), (b) *Max*(*Tcp*).

erage \triangleright Both average delay and maximum delay of SDP achieve the similar performance.

200 co **Run-time comparison**

Self-adaptive Quadruple Partition

 (a) (b) Figure: Run time comparison among different strategies, (a) with ILP, (b) with TILA.

 \triangleright The run-time of SDP is lower than ILP but slower than TILA [2] (Yu. et al. ICCAD'2015).

Figure: Sub-grid partition illustration, (a) sub-grid partition, (b) sub-grid partition tree.

Semidefinite Programming Relaxation

The proposed self-adaptive partition provides an opportunity for further speed-up. Semidefinite programming (SDP) is solvable in polynomial time while providing a theoretically better solution than Linear Programming (LP).

 \triangleright The objective function:

16 de janeiro de la construction de la construction de la construction de la construction de la construction
16 de janeiro de la construction de la construction de la construction de la construction de la construction
16 d This work is supported in part by NSF, Oracle, and CUHK Direct Grant for Research.

 $\bigg)$

 $min(T \cdot X)$. ▶ Matrix *T* - $|S \cdot L|$ -dimension symmetric matrix representing timing costs

 \triangleright Matrix $X - |S \cdot L|$ -dimension symmetric matrix representing variables

 $t_s(i, j)$... $t_v(i, j, p, q)$

 $\bigg)$

 X_{ij} ... Y_{ijpq}

Semidefinite Programming Example

Figure: Example of layer assignment through solving SDP.

 \triangleright Cost matrix T and solution matrix X of this example

■ Segment S2 overlaps with other segments, resulting in continuous solutions.

 \triangleright Therefore, post mapping is required to provide integer solutions.

Overall Algorithm Flow

Post mapping transfers continuous solutions into discrete assignment result.

Mapping Flow

Figure: Iterative incremental layer assignment algorithm flow.

Experimental Results

Experimental Setup

- ► CPLA implemented in C++, Gurobi as the MILP solver and CSDP as the SDP solver
- \blacktriangleright Linux machine with 2.9GHz Inter(R) Core and 192GB memory
- Initial routing and layer assignment result from NCTU-GR and NVM tool

Comparison on layer assignment result with TILA-0.5%

Figure: Comparison among different strategies, (a) *Avg*(*Tcp*), (b) *Max*(*Tcp*), (c) via number, (d) via violations.

 \triangleright Average delay improved by 14% while maximum delay improved by 4%. \triangleright Via violations improved by 10% with the similar number of vias.

Performance comparison with ILP

Partition size impact on performance

Segment# in each partition

 \triangleright When partition size is equal to 10, the lowest run-time can be obtained.

Figure: Comparison among different partition sizes, (a) *Avg*(*Tcp*), (b) *Max*(*Tcp*), (c) run time.

► For different partition sizes, similar performance can be achieved.

Acknowlegement