

Hierarchical Whitespace Allocation in Top-down Placement*

Andrew E. Caldwell, Andrew B. Kahng and Igor L. Markov

Cadence Design Systems, River Oaks Parkway, San Jose, CA 95134

UC San Diego CSE and ECE Departments, La Jolla, CA 92093-0114

Univ. of Michigan, EECS Dept., Ann Arbor, MI 48109-2122

caldwell@cadence.com abk@ucsd.edu imarkov@umich.edu

Abstract

Increased transistor density in modern commercial ICs typically originates in new manufacturing and defect prevention technologies [15, 16]. Additionally, better utilization of such *low-level* transistor density may result from improved software that makes fewer assumptions about physical layout in order to reliably automate the design process. In particular, recent layouts tend to have large amounts of whitespace, which is not handled properly by older tools. We observe that a major computational difficulty arises in partitioning-driven top-down placement when regions of a chip lack whitespace. This tightens balance constraints for min-cut partitioning and hampers move-based local-search heuristics such as Fiduccia-Mattheyses. However, the local lack of whitespace is often caused by very unbalanced distribution of whitespace during previous partitioning, and this concern is emphasized in chips with large overall whitespace.

This paper focuses on accurate computation of tolerances to ensure smooth operation of common move-based iterative partitioners, while avoiding cell overlaps. We propose a mathematical model of hierarchical whitespace allocation in placement, which results in a simple computation of partitioning tolerance purely from relative whitespace in the block and the number of rows in the block. Partitioning tolerance slowly increases as the placer descends to lower levels, and relative whitespace in all blocks is limited from below (unless partitioners return “illegal” solutions), thus preventing cell overlaps. This facilitates good use of whitespace when it is scarce and prevents very dense regions when large amounts of whitespace are available.

Our approach improves the use of the available whitespace during global placement, thus leading to smaller whitespace requirements. Existing techniques, particularly those based on Simulated Annealing [21, 10], can be applied after global placement to bias whitespace with respect to particular concerns, such as routing congestion, heat dissipation, cross-talk noise and DSM yield improvement.

1 Introduction

The progression of Moore’s law [18, 15] for commercial ICs has been so far maintained by steady increases in device densities as a result of innovations in manufacturing and defect prevention technologies [16]. At the same time, device density for a given process generation is also limited by the capabilities of EDA software and the assumptions made by software developers.

*This work was performed when the authors were affiliated with UCLA. All implementations are publicly available from <http://vlsicad.cs.ucla.edu/software/PDtools/>

Historically, utilization rates increased (i.e., whitespace decreased) steadily with the introduction of 3-, 4- and even 5-layer metal technologies. In contrast with the preceding 2-layer metal regime, 3 or more layers of metal brought the following changes: (1) the need for routing channels disappeared, (2) double-back (shared power and ground rail) standard-cell styles removed all whitespace between cell rows, and (3) built-in cell library porosity became the proxy for previous techniques such as explicit feedthrough insertion. By the late 1990's, standard-cell P&R blocks exhibited a clear trend of decreasing whitespace, and instances with single-digit percentages of whitespace were not uncommon. On the other hand, as CMOS technology moved below quarter-micron feature sizes, utilization rates started to decrease (i.e., whitespace has *increased*) in 180nm and 130nm designs. One obvious reason for the increase in whitespace is the increasingly interconnect-limited nature of designs.

Interconnect limits arise due to pad-limited designs, and also due to conflicting goals (power and clock distribution, reliability, signal density, etc.) for the UDMSM interconnect architecture. (For example, there is little point in pushing transistor density to its limits if the resulting cell library has inadequate porosity, or if power/ground distribution then uses up too much of the local routing resources.) Another reason for increased whitespace stems from the growing number of macro blocks in SoC designs which, even when floorplanned well, leave large layout areas in which to place the remaining standard-cell logic (e.g., "block limited" designs). However, increased whitespace is also a consequence of limitations in place-and-route tools. Guardbanding for these limitations — in particular, the ability to deal with routing congestion — in order to minimize the number of passes through the back-end tool flow is a major contributor to whitespace at the 180nm and 130nm nodes. Recent anecdotal evidence indicates that whitespace today varies considerably (from 20+% to around 70+%) with methodology and allowed design time.

Global top-down partitioning-based placers can be sensitive to the amount of whitespace, which contributes to tool limitations. As the *partitioning tolerance* decreases in the partitioning instance, the results of modern hypergraph partitioners (cutsizes, and the ability to find a solution that meets balance constraints) deteriorate significantly [8]. Partitioning solutions that violate prescribed tolerance often lead to cell overlaps when future subproblems run out of whitespace. Since non-uniform cell sizes generally worsen partitioner performance, modern cell libraries that have widely varying cell sizes exacerbate the difficulties of move-based partitioners. The role of partitioners in limiting the reduction of whitespace is the focus of our present research.

The essential components of a typical placement problem are the *placement region*, possibly with discrete allowed locations, the *cells* to be placed subject to various constraints, and the *netlist topology* that shapes the minimized objective function. Academic and commercial standard-cell placers often apply a top-down, divide-and-conquer approach to define an initial *global placement*. The top-down approach [20, 13, 2] decomposes the given placement problem into smaller problems by subdividing the placement region, assigning cells to subregions, reformulating constraints, and cutting the netlist — such that good solutions to subproblems combine into good solutions of the original problem. The concept of *placement blocks* is pivotal. A block represents (i) a placement region with allowed locations, (ii) a collection of cells to be placed in this region, (iii) all nets incident to the cells, and (iv) locations of all cells beyond the given region that are adjacent to the cells to be placed in the region; such external cells are considered to be terminals for the block, and their locations are fixed. A high-level pseudocode for top-down global placement in terms of

placement blocks is shown in Figure 1; sufficiently small partitioning instances are processed as *end-cases* by specialized partitioners or placers.

Every block yields a hypergraph partitioning instance: nodes correspond to cells inside the block as well as propagated external terminals [7], and hyperedges are induced over the node set from the original netlist. In practice blocks are split through balanced min cut hypergraph bisection with FM-type move-based heuristics [12, 9]; The performance of such heuristics on larger instances is improved through the multilevel paradigm [3, 11]. Global placement solutions place all cells near legal sites in cell rows with minimal overlaps. Detailed placement refinement then makes small perturbations to the placement to legalize (remove overlaps) or improve routability. The more overlapping cells and unroutable hotspots the detailed placer must address, the larger the risk of solution quality degradation. As more constraints or optimization objectives are added, e.g., timing, density or IR-drop, legalizing an overlapping or unroutable global placement while preserving solution quality becomes increasingly difficult. In many cases, the problem may not be fixable at the detailed placement stage. To this end, it is reasonable to aim for non-overlapping, routable global placements. One common technique for managing both legality and routability during placement is controlling the distribution of whitespace – inserting more whitespace in congested areas. In this work, we focus on mathematical methods of whitespace management during top-down placement with the goal of producing non-overlapping, routable global placements.

Hypergraph bisection instances arising in top-down the placement process in some cases have tight *balance constraints* [8], i.e., the sizes of partitions in the solution should not deviate from *target* partition sizes by more than prescribed amounts.¹ Such constraints may arise in several circumstances. In leading-edge microprocessors, the proportion of deliberately introduced free sites, i.e., whitespace, is limited.² In other designs whitespace may be distributed unevenly, and thus some regions may be quite dense. To avoid overcapacity blocks in such regions, total cell area assigned to a block must closely match the area available for cells. Indeed, illegal solutions or excessively relaxed balance tolerances lead to uneven area utilization (i.e., proportion of whitespace in a given sub-block) and, eventually, overlapping cells. Even when a legal solution is obtained by the partitioner, as the partitioner exploits its available tolerance one child of the partitioned block can have relatively less whitespace than its parent, so that partitioning constraints become tighter on lower levels of top-down placement *even if average whitespace is large*. Such deviations in available whitespace cannot be easily corrected, essentially because cuts parallel to rows cannot be adjusted after partitioning. On the other hand, attempting to maintain available whitespace by imposing unnecessarily tight balance constraints will hurt solution quality and lead to increased wirelength in the layout.³

In this paper we develop a mathematical model to find better tolerances: not too big — to generally avoid overcapacity or unroutable blocks, and not too small — to facilitate common move-based partitioners. *The*

¹See [1] for a review of netlist partitioning formulations and constraints.

²Deliberate introduction of whitespace may be related to pin-limited designs, limits of power distribution/dissipation density, and the need to maintain autoroutability. While these phenomena are increasingly addressed by packaging and process technology, better architecture and circuit design techniques, multilayer interconnect processes, and cell library design, block-limited designs are becoming increasingly common. This accounts for the large range of whitespace currently seen in SoC designs.

³The work of [8] in particular showed that iterative move-based partitioners perform poorly with small tolerance. The authors of [8] proposed the technique of *intermediate relaxations* to trade CPU time for solution quality in such circumstances. Their work is orthogonal to ours: we address the question of how to best control the allocation of whitespace (via bounds on partition sizes) during the top-down partitioning process.

```

Variables:      A queue of blocks
Initialization: A single block represents the original placement problem
Algorithm:     while (queue not empty)
                dequeue a block
                if (small enough) consider endcase
                else
                    bipartition into smaller blocks
                    enqueue each block

```

Figure 1: High-level outline of the top-down partitioning-based placement process.

primary contribution of this paper is a new whitespace management framework for global placement whose utility is supported by experimental evidence. The framework can include consideration of layout density issues, such as routability and IR-drop, during the top-down placement process.

Whitespace allocation can be further improved using recent annealer-based algorithms for detailed placement that address routing congestion [21, 10], heat dissipation, cross-talk noise, DSM yield improvement, etc. Modifications to global placement have been proposed targeting wirelength [4] and congestion [17]. We believe that our mostly mathematical contribution, being an easy addition to top-down partitioning-based placers, provides a common denominator for such work.

The remaining text is organized as follows. Basic notation is introduced in Section 2 and straightforward facts about whitespace are mentioned. In Section 3, we control splitting of single blocks by maximizing partitioning tolerance for given *whitespace deterioration*. This process is described mathematically. Most surprisingly, *relative whitespace* and *relative tolerances* are connected independently of whitespace deterioration, site and cell areas. This leads to a simple computation of bipartitioner parameters in terms of relative whitespace in the block and the number of rows in the block. Experimental validation is presented in section 4. Section 5 concludes with directions for ongoing work.

2 Whitespace fundamentals

Let the block have *site area* S , cell area C (typically with $C \leq S$), *absolute whitespace* $W = \max\{S - C, 0\}$ and *relative whitespace* $w = W/S$.⁴ A geometric bipartitioning of the block entails *site areas* S_0 and S_1 in child blocks, such that $S_0 + S_1 = S$, $S_0 \leq S$ and $S_1 \leq S$. Any hypergraph bipartitioning solution implies *cell areas* C_0 and C_1 in child blocks, such that $C_0 + C_1 = C$, $0 \leq C_0$, $0 \leq C_1$ (if in the original block $C \leq S$, we will also require $C_0 \leq S_0$, $C_1 \leq S_1$). The input to a hypergraph bipartitioner must specify both the netlist and the allowed ranges for C_0 and C_1 , i.e. bounds $C_0^{min} \leq C_0 \leq C_0^{max}$, $C_1^{min} \leq C_1 \leq C_1^{max}$ (with $C_0^{max} + C_1^{max} > C$ possible). These bounds establish *absolute* $T_j = C_j^{max} - C_j^{min}$ and *relative* $\tau_j = T_j/C$ partitioning tolerances for $j = 0, 1$

($\tau_0 = \tau_1$ not required, but often holds). *Absolute* and *relative* whitespace in child blocks are defined by $W_j = S_j - C_j$ and $w_j = \frac{W_j}{S_j}$ (see Figure 2).

⁴We will sometimes refer to W as *whitespace*.

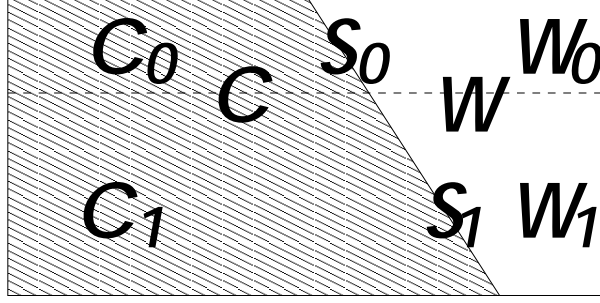


Figure 2: Basic variables for a block and two child blocks: site area S , cell area C and whitespace $W = S - C$.

In order to express w in terms of w_0 and w_1 , write $w = \frac{S-C}{S} = \frac{S_0-C_0}{S} + \frac{S_1-C_1}{S} = \frac{S_0}{S} \frac{W_0}{S_0} + \frac{S_1}{S} \frac{W_1}{S_1}$, hence

$$w = \frac{S_0}{S} w_0 + \frac{S_1}{S} w_1 \quad (1)$$

Since $\frac{S_1}{S} = 1 - \frac{S_0}{S}$, relative whitespace in a block is a convex combination of relative whitespace in child blocks. Subsequently relative whitespace in a block is never smaller than that in every child block. If one child block has more relative whitespace than its parent, then the other child block has less. If we wish to limit relative whitespace in blocks from below, it suffices to consider end-case blocks only. In the special case when child blocks have equal site area, Equation (1) contains an arithmetic average. In practice S_0/S can be small, e.g., $1/3$ when a three-row block is bipartitioned horizontally, and even smaller than $1/3$ when some sites in the one-row child block are obstructed by special wiring.

Given a collection of non-overlapping placement blocks that cover the layout, e.g., in the course of top-down placement, one can recursively apply (1), to show that the *average relative whitespace* for the design (i.e. relative whitespace for the top-level block) is a convex combination of relative whitespace in individual blocks.

Overcapacity cell area (or *overflow*) in a block is defined by $\Theta = \max\{0, C - S\}$, *relative overflow* — by $\theta = \Theta/S$. Given a collection of non-overlapping blocks that cover the layout, the *total overflow* Θ^0 for the design is the sum of overflow in all blocks; *average relative overflow* θ^0 for the design is computed by dividing Θ^0 by the site area in all blocks. Similarly to average relative whitespace, average relative overflow is a convex combination of relative overfills in individual blocks.

For a given block, the relative whitespace and relative overflow can not be simultaneously bigger than zero, and ensuring non-zero whitespace in all blocks precludes overflow. Assuming non-zero relative whitespace [at the top-level], we will require that for each block split the relative whitespace in each child block is at least αw , where w is the relative whitespace in the block and $0 \leq \alpha \leq 1$ is *whitespace deterioration*, i.e.

$$w_0 \geq \alpha w \text{ and } w_1 \geq \alpha w \quad (2)$$

For practical purposes, $\alpha = 1$ and may be overly restrictive as it entails partitioning with zero tolerance in the proportion of site area, while $\alpha = 0$ may be too loose as it allows for child blocks with no whitespace. As α approaches one the distribution of whitespace in the final placement approaches uniform. An α of zero

allows for fully utilized regions of the layout. This can improve wirelength and timing, but may also result in routability problems. It is desirable, then, to adjust α on a per-block basis to account for maximum allowed layout densities in leaf-level blocks. For example, in areas of high predicted routing density α should be high, whereas in sparser areas it can safely be set to a low value α allowing for denser placement but still ensuring a non-overlapping result.

The following is the key observation for hierarchical whitespace allocation in top-down placement.

Theorem Assuming non-zero average relative whitespace w_0 in the design, the result of top-down placement will have zero total overfill if every block split is performed with whitespace deterioration $\alpha_i > 0$, possibly different for every block.

Proof Relative whitespace in every end-case block will then be at least $\alpha_0 \dots \alpha_n w_0 > 0$.

In partitioning-driven top-down placement, whitespace deterioration is controlled with balance tolerance that constrains partitioning solutions.

Corollary Top-down placement results in a zero-overfill placement if balance tolerances correspond to strictly positive whitespace deterioration and all partitioning solutions are legal.

Excessively tight (i.e. to close to 1) whitespace deterioration may allow no legal solutions, e.g., for purely number partitioning reasons if it entails tolerance below the size of one site.⁵ Small tolerance can also imply poorer partitioning quality because it restricts the solution space and may incapacitate move-based partitioners by disallowing moves of large cells (when their sizes are bigger than tolerance).

On the other hand, [5] empirically shows that higher partitioning tolerances result in small placement wirelength. Therefore it is important to keep whitespace high enough to prevent overcapacity blocks and, at the same time, keep partitioning tolerances high to ensure small cuts and, thus, wirelength.

3 Whitespace allocation

In this section, we will control block splits by maximizing partitioning tolerance for given whitespace deterioration; this results in a number of useful properties: (6), (7) and (8). Most surprisingly, relative whitespace and relative tolerances are recursively connected (12) independent of whitespace deterioration. This allows to compute optimal relative tolerance given the initial and final relative whitespace. The latter can be computed solely from the initial relative whitespace (10).

We show how to determine C_j^{min} and C_j^{max} for a block

only from its relative whitespace and the number of rows covered (16). Our formulae, derived to account for the worst case possible, can be transparently adjusted for optimism, which may be useful in designs with abundant whitespace.

3.1 Splitting a block

We compute C_j^{max} and C_j^{min} for child blocks of a particular block assuming given α . The resulting tolerances for both partitions appear equal and thus the relative tolerance τ is determined by α . It is somewhat surprising that α , w and τ are connected independently of S_j or C_j and any two of them imply a particular value for the

⁵Cell sizes are typically small integer multiples of site size.

third, which can be computed by a simple formula. These relations facilitate further modeling of recursive block splits in the next subsection.

Rewrite (2) as

$$\alpha \frac{S-C}{S} \leq \frac{S_j-C_j}{S_j} \Rightarrow C_j S \leq (1-\alpha)S_j S + \alpha S_j C \Rightarrow C_j \leq (1-\alpha)S_j + \alpha \frac{S_j}{S} C$$

adding the original bounds for C_j we get

$$0 \leq C_j \leq \min\{C, S_j, (1-\alpha)S_j + \alpha \frac{C}{S} S_j\}$$

Since $0 \leq \alpha \leq 1$ and $0 \leq C \leq S$, we have

$$(1-\alpha)S_j + \alpha \frac{C}{S} S_j \leq (1-\alpha)S_j + \alpha S_j \leq S_j.$$

The above is now simplified⁶

$$0 \leq C_0 \leq \min\{C, (1-\alpha)S_0 + \alpha \frac{C}{S} S_0\} =: C_0^{max} \quad (3)$$

$$0 \leq C_1 \leq \min\{C, (1-\alpha)S_1 + \alpha \frac{C}{S} S_1\} =: C_1^{max} \quad (4)$$

The remaining constraint $C_0 + C_1 = C$ is now equivalent to

$$C_0 \geq \max\{0, C - C_1^{max}\} =: C_0^{min}, \quad C_1 \geq \max\{0, C - C_0^{max}\} =: C_1^{min} \quad (5)$$

When C is very small compared to S (i.e. the block has a lot of whitespace) and α sufficiently small, C_j^{max} and C_j^{min} may degenerate into C and 0 respectively. In such cases, all cells are allowed to go into one partition and no further analysis is required until a child block appears with small enough white space to produce non-trivial partitioning tolerance.

In the following analysis, we assume that all cells can never go into one partition (worst case), therefore⁷ $C_j^{max} = (1-\alpha)S_j + \alpha \frac{S_j}{S} C$ and $C_j^{min} = C - C_{1-j}^{max}$. Now $T_j = C_j^{max} - C_j^{min} = C_0^{max} + C_1^{max} - C = (1-\alpha)(S-C)$, i.e. absolute tolerances for partitions are equal. Furthermore, $\tau = T/C = (1-\alpha)(S/C - 1)$ from which straightforward calculations lead to

$$\tau = \frac{(1-\alpha)w}{1-w} \quad (6)$$

$$\alpha = (\tau + 1) - \frac{\tau}{w} \quad (7)$$

$$w = \frac{\tau}{\tau + 1 - \alpha} \quad (8)$$

3.2 Hierarchical whitespace allocation

It has been shown above that, for a given block, feasible ranges for partition capacities are uniquely determined by α . This section discusses methods of determining values of α . The methods presented do not account for routability, as this will depend heavily on the capabilities of the specific router to be used. The

⁶Here we also define C_0^{max} and C_1^{max} , while C_0^{min} and C_1^{min} are defined in (5).

⁷Our analyses hold for blocks allowing all cells in one partition, however, the proposed values of C_j^{max} on lower levels are not the best possible. In a way, we assume the pessimistic case.

framework presented can be, however, simply extended to utilize congestion estimation by increasing or decreasing the target whitespace of the leaf-level blocks or the optimism applied to α . Further, we assume layout cut-lines are shifted to match w_0 and w_1 as closely as possible to w in an attempt to recover the whitespace degradation. This allows future partitioning steps the benefits of high-tolerance, but may be inappropriate for low- utilization designs as it attempts to distribute whitespace uniformly resulting in an overly-spread placement.

We start with the top-level block having w_0 whitespace and go on to further blocks with whitespace

$$w_{i+1} = \alpha_i w_i \quad (9)$$

We distinguish between blocks being split by cut lines parallel to rows and those perpendicular to rows. This is because perpendicular cut lines can be adjusted after partitioning to achieve nearly any desired distribution of whitespace among child blocks, which has the effect of $\alpha \approx 1$. In other words, we only have to consider row-parallel cut-lines. Let R be the number of rows; the expected number n of *recursively applied parallel block splits* will be $n = \lceil \log_2 R \rceil$, assuming that rows are distributed evenly between child blocks at every block split.

To prevent overcapacity blocks and improve routability, assume also that whitespace in every block below the current block must be at least \bar{w} ($\leq w_0$). We find \bar{w} , observing that end-case blocks have $\alpha = 0$ since they are not partitioned further. Thus from (8) we get

$$\bar{w} = w_n = \frac{\bar{\tau}}{\bar{\tau} + 1} \quad (10)$$

Note that \bar{w} can be determined individually for each block being partitioned, and provides a method of accounting for routability and similar density-related issues. A straightforward way to determine α_i is to assume that they are all equal. This leads to⁸

$$\alpha = \sqrt[n]{\frac{\bar{w}}{w_0}} = \sqrt[n]{\frac{\bar{\tau}}{w_0(\bar{\tau} + 1)}} \quad (11)$$

However, assuming all τ_i equal appears more practical since balanced partitioners typically require a certain relative tolerance to be successful regardless of whitespace deterioration. In a different, improved approach we combine (7) and (9) to get rid of α

$$w_{i+1} = (\tau_i + 1)w_i - \tau_i \quad (12)$$

Now, assuming that all τ_i equal yields $w_n = (\tau + 1)^n w_0 - \tau(\tau + 1)^{n-1} - \dots - \tau = (\tau + 1)^n w_0 - \tau \frac{(\tau + 1)^n - 1}{(\tau + 1) - 1}$ resulting in

$$\bar{w} = w_n = (\tau + 1)^n w_0 - (\tau + 1)^n + 1$$

and

$$(\tau + 1)^n = \frac{1 - \bar{w}}{1 - w_0} \quad (13)$$

⁸Due to essentially uncontrollable distribution of whitespace in child blocks and to compensate for small whitespace fluctuations after perpendicular block splits, α_i can be recomputed for every block, given the actual amount of whitespace in the block.

Therefore we can replace the straightforward Equation (11) with

$$\tau = \sqrt[n]{\frac{1-w}{1-w_0}} - 1 \quad (14)$$

and combine (13) and (10) with $\tau_n = \tau$ to find a closed expression for τ

$$\tau = \frac{1}{\sqrt[n+1]{1-w_0}} - 1 \quad (15)$$

Finally, (15) and (7) give a closed expression for whitespace deterioration α in terms of relative whitespace w in the current block and the number R of rows in the block.

$$\alpha = \frac{\sqrt[n+1]{1-w} - (1-w)}{w \sqrt[n+1]{1-w}}, \quad n = \lceil \log_2 R \rceil \quad (16)$$

Values C_j^{max} and C_j^{min} ($j = 0, 1$) supplied to partitioner can be directly computed by (3), (4) and (5).

This computation of C_j^{max} and C_j^{min} can be performed for every block using the exact amount of whitespace available after partitioning.⁹ While resulting tolerances are likely to increase towards lower levels, they will preserve original lower bounds for relative whitespace in every block¹⁰ and thus prevent overcapacity blocks. Compared to the pessimistic *a priori* tolerances, such increased tolerances can result in cut improvements during hypergraph partitioning and cause smaller total placement wirelength.

3.3 Splitting overcapacity blocks

While our approach guarantees no overcapacity blocks and no cell overlaps under certain conditions, these conditions may not always hold. In particular, a partitioner may return illegal solutions when there are cells larger than partitioning tolerance or, more generally, when balanced solutions do not exist for number partitioning reasons. This typically happens when partitioning small blocks with tight tolerances, where each cell accounts for a considerable percent of the block's total cell area.

When splitting such blocks, we will minimize the maximal relative overfill in child blocks, which is equivalent to equal relative overfill according to Section 2.¹¹ Hence, the desired partition capacities C_1 and C_2 need to satisfy $\theta^0 = \frac{C_0 - S_0}{S_0} = \frac{C_1 - S_1}{S_1} = \theta^1$ and $C_0 + C_1 = C$; they can be computed via

$$C_0 = \frac{S_0}{S} C, \quad C_1 = \frac{S_1}{S} C \quad (17)$$

These considerations do not allow for nonzero partitioning tolerance since even the slightest imbalance in the resulting partitioning solution would lead to an increase in relative overfill. However, there are currently no practical algorithms available for zero-tolerance partitioning. Therefore, it is necessary to artificially introduce partitioning tolerance. Appropriate values should be chosen corresponding to the cell area distribution and the capabilities of the partitioning algorithm used, e.g., we used the smaller of 2% and the overall whitespace in the design. Formulae (17) are used for computing target partitioning capacities.

⁹Compared to a possible *a priori* computation for all blocks that only assumes relative whitespace at the top level and the number of rows in the design, but assumes that *every block has the worst possible relative whitespace after partitioning*.

¹⁰In the assumption that all partitioning tolerances are satisfied.

¹¹Since the relative overfill in the original block is a convex combination of the relative overfills in child blocks, one of the child blocks having smaller overfill implies the other having a higher overfill.

4 Experimental validation

In this section we present a comparison of our proposed method with a simpler alternative whitespace allocation strategy for horizontal cuts. In this simpler method, the target partition balances are proportional to site areas in the resulting partitions and horizontal cut tolerances are constant.

4.1 Top-down placement testbench

We have implemented a full-fledged global placer that reads standard-cell row-based designs from Cadence Design Systems, Inc. in LEF/DEF formats and produces cell placements with little or no overlap. Placement blocks are split via min-cut partitioning, which is implemented as a variant of multilevel Fiduccia-Mattheyses [9] for blocks with 200 cells or more and [flat] Fiduccia-Mattheyses for smaller blocks. The placer chooses vertical or horizontal block splits depending on the blocks' aspect ratio to always cut along the longest side of a block. When partitioning is performed with vertical cut line, the current block is bisected by a straight line and sites in the two resulting regions are counted to produce the site area in each region. The target partition capacities (cell areas) are then computed to be proportional to the site areas and add up to the cell area in the block. Vertical partitioning is performed with 10% tolerance in all our experiments. After partitioning, when the actual total cell area in each partition is available, the vertical straight line determining block boundaries is optimally shifted to equalize relative whitespace in the blocks.

We use two different methods to allocate whitespace for blocks split by horizontal cut lines. In one set of experiments we simply used the same computations as for the vertical block splits with partitioning tolerance being 2%, except for the inability to adjust horizontal cut lines after partitioning.¹² The second method follows Formulae (16), (3), (4) and (5). All overcapacity blocks are treated as explained in Subsection 3.3.

We do not apply “legalization”, “cyclong”, “overlapping” [14] or any other techniques that would move cells between existing blocks or change block boundaries other than during top-down block splits. Without such post-processing the effects of different approaches to whitespace allocation come clear, and fair comparisons can be made. Moreover, even with such pessimistic view of cell overlaps, our top-down placements have very few cell overlaps. Our top-down placer is implemented in the C++ language. Executables are compiled with the SunPro CC4.2 compiler on Solaris2.6 and optimized with -O5. We use Sun Ultra-10 300MHz workstations with 256Mb of memory.

| Test Case | Core cells | Pads | Nets | core rows | % whitespace |
|-----------|------------|------|--------|-----------|--------------|
| case1 | 2741 | 545 | 2842 | 24 | 11.29 |
| case2 | 11471 | 662 | 11673 | 42 | 24.28 |
| case3 | 12146 | 711 | 10880 | 53 | 23.40 |
| case4 | 20392 | 185 | 21987 | 88 | 14.07 |
| case5 | 85395 | 177 | 87272 | 301 | 30.15 |
| case6 | 117440 | 177 | 101726 | 250 | 14.21 |

Table 1: Statistics for the testcases used in our experiments.

¹²Which explains our choice of smaller tolerance.

4.2 Results

We report experimental results for six industrial testcases, ranging from $2.7Kcells$ to more than $117Kcells$, with whitespace from 11% to 30% (see Table 1).¹³ Relative average overfill for final placements is evaluated as explained in Section 2 and is given together with final wirelength and CPU time in Table 2.

| | Proposed | | | Fixed Tolerance 5% | | | Fixed Tolerance 2% | | |
|-------|----------|-------|------|--------------------|-------|------|--------------------|-------|------|
| | HPWL | CPU | R/O | HPWL | CPU | R/O | HPWL | CPU | R/O |
| case1 | 5.85e5 | 28.5 | 0.00 | 5.95e5 | 27.4 | 0.00 | 6.32e6 | 28.1 | 0.00 |
| case2 | 2.75e5 | 158.2 | 0.00 | 2.80e5 | 151.1 | 0.00 | 3.03e5 | 158.7 | 0.00 |
| case3 | 2.50e5 | 140.5 | 0.00 | 2.61e5 | 140.3 | 0.00 | 2.89e5 | 153.8 | 0.00 |
| case4 | 5.83e6 | 298.3 | 0.03 | 5.83e6 | 282.9 | 0.04 | 6.38e6 | 295.1 | 0.00 |
| case5 | 2.26e7 | 1894 | 0.00 | 2.29e7 | 1802 | 0.00 | 2.45e7 | 1784 | 0.00 |
| case6 | 1.16e8 | 2479 | 0.01 | 1.16e8 | 2438 | 0.08 | 1.25e8 | 2424 | 0.00 |

Table 2: Placer run-time (CPU seconds on a 300MHz Sun Ultra-10), half-perimeter wirelength of final placements and average relative overfill (**R/O**) in percent for our proposed tolerance computation method and a straightforward method with fixed tolerance. All numbers are averages of five runs.

Our experiments show that

- Our proposed method allows for maximal use of available whitespace, producing placements with typically 10% smaller wirelength than the fixed 2% method.
- Increasing the fixed tolerance enough to achieve competitive wirelengths produces placements with more overlap than our proposed method.
- For every test case, our proposed method produces either lower wirelength and the same amount of relative overlap, or no relative overlap and lower wirelength, than the fixed 5% method.

5 Conclusions and future work

Our work addresses the global placement of modern ICs in which whitespace can vary dramatically across the layout region. The proposed methods allow handling very small amounts of whitespace and can be used to appropriately distribute large amounts of whitespace when it is available. Our framework can be extended to account for target utilizations at particular areas of the layout region, allowing for consideration of density-related effects such as routability and IR-drop.

We have derived simple formulae for optimal¹⁴ hierarchical whitespace allocation in top-down partitioning-driven placement of standard-cell row-based ASIC designs. With a straightforward practical adjustment, partitioning tolerance slowly increases as the placer descends to lower levels. We limit relative whitespace in all blocks from below, and this constraint prevents overcapacity blocks under the assumption that all partitioning solutions are legal. We point out that *using the proposed formula requires almost no additional*

¹³White space is measured by dividing the site area available to cells by the total cell area in the netlist. In this we accounted for sites put out of use by power stripes and other obstacles, but such adjustments turn out to be very small.

¹⁴In the worst-case sense; also under the assumption of constant partitioning tolerance and legal partitioning solutions on all levels.

programming and incurs practically no runtime penalty. The use of this formula also does not appear to conflict with popular placement algorithms and design practices.

Experiments on industrial testcases with up to 114K cells show that our technique practically achieves cell overlaps on the order of hundredths of a percent of the total areas and is superior to a commonly used straightforward technique since it achieves better wirelength under the assumption of comparable cell overlaps (which can be provided by sufficiently small partitioning tolerance).

In real-world implementations, whitespace allocation can be further improved using recent annealer-based algorithms for detailed placement that address routing congestion [21, 10], heat dissipation or other concerns. Modifications to global placement have been proposed targeting wirelength [4] and congestion [17]. We believe that our mostly mathematical contribution, being an easy addition to top-down partitioning-based placers, provides a common denominator for such work without any noticeable costs. This can be contrasted with non-trivial development required to implement algorithms based on simulated annealing and their runtime costs [21].

References

- [1] C. J. Alpert and A. B. Kahng, "Recent Directions in Netlist Partitioning: A Survey", *Integration*, 19(1995) 1-81.
- [2] C. J. Alpert, T. Chan, D. J.-H. Huang, I. Markov and K. Yan, "Quadratic Placement Revisited", *Proc. ACM/IEEE Design Automation Conference*, 1997, pp. 752-757.
- [3] C. J. Alpert, J.-H. Huang and A. B. Kahng, "Multilevel Circuit Partitioning", *ACM/IEEE Design Automation Conference*, 1997, pp. 530-533.
- [4] C. J. Alpert, G.-J. Nam and P. G. Villarrubia, "The Effect of Free Space in Global Placement", to appear in *ICCAD 2002*.
- [5] A. E. Caldwell, A. B. Kahng and I. L. Markov, "Relaxed Partitioning Balance Constraints in Top-Down Placement", *In Proc. ASIC '98*
- [6] W. Deng, personal communication.
- [7] A. E. Dunlop and B. W. Kernighan, "A Procedure for Placement of Standard Cell VLSI Circuits", *IEEE Transactions on Computer-Aided Design* 4(1) (1985), pp. 92-98
- [8] S. Dutt and H. Theny, "Partitioning Around Roadblocks: Tackling Constraints With Intermediate Relaxations", *Proc. IEEE International Conference on Computer-Aided Design*, 1997, pp. 350-355.
- [9] C. M. Fiduccia and R. M. Mattheyses, "A Linear Time Heuristic for Improving Network Partitions", *Proc. ACM/IEEE Design Automation Conference*, 1982, pp. 175-181.
- [10] B. Hu and M. Marek-Sadowska, "Congestion Minimization During Placement Without Estimation", to appear in *ICCAD 2002*.
- [11] G. Karypis, R. Aggarwal, V. Kumar, and S. Shekhar, "Multilevel Hypergraph Partitioning: Applications in VLSI Design", *Proc. ACM/IEEE Design Automation Conference*, 1997, pp. 526-529.
- [12] B. W. Kernighan and S. Lin, "An Efficient Heuristic Procedure for Partitioning Graphs", *Bell System Tech. Journal* 49 (1970), pp. 291-307.
- [13] J. Kleinmans, G. Sigl, F. Johannes, K. Antreich, "GORDIAN: VLSI Placement by Quadratic Programming and Slicing Optimization", *IEEE Trans. on Computer Aided Design* 10(3) (1991), pp. 356-365.
- [14] D. J. Huang and A. B. Kahng, "Partitioning-Based Standard Cell Global Placement with an Exact Objective", *Proc. ACM/IEEE International Symposium on Physical Design*, 1997, pp. 18-25.
- [15] David Jensen, Charles Gross, Dinesh Mehta, "New industry document explores defect reduction technology challenges", *Micro Magazine*, January 1998, <http://www.micromagazine.com/archive/98/01/jensen.html>
- [16] David Jensen and William Fosnight, "Defect prevention and elimination: Where the rubber hits the road(map)", *Micro Magazine*, October 1998, <http://www.micromagazine.com/archive/98/10/jensen.html>
- [17] A. Rohe, U. Brenner, "An Effective Congestion Driven Placement Framework", *ISPD 2002*, pp. 6-11.
- [18] Semiconductor Industry Association, *National Technology Roadmap for Semiconductors*, San Jose, CA, SIA, 1997.
- [19] H. D. Simon and S.-H. Teng, "How Good is Recursive Bisection?", *SIAM J. Scientific Computing* 18(5) (1997), pp. 1436-1445.
- [20] R. S. Tsay and E. Kuh, "A Unified Approach to Partitioning and Placement", *IEEE Trans. on Circuits and Systems*, 38(5) (1991), pp. 521-633.
- [21] X. Yang, B.-K. Choi and M. Sarrafzadeh, "Routability-driven Whitespace Allocation for Fixed-Die Standard Cell Placement". *ISPD 2002*, pp. 42-48.