



Master-Thesis

High Quality Hypergraph Partitioning via Max-Flow-Min-Cut Computations

Tobias Heuer

Advisors: Prof. Dr. Peter Sanders
M. Sc. Sebastian Schlag

Institute of Theoretical Informatics, Algorithmics
Department of Informatics
Karlsruhe Institute of Technology

Date of submission: 22.12.2017

Hiermit versichere ich, dass ich diese Arbeit selbstständig verfasst und keine anderen, als die angegebenen Quellen und Hilfsmittel benutzt, die wörtlich oder inhaltlich übernommenen Stellen als solche kenntlich gemacht und die Satzung des Karlsruher Instituts für Technologie zur Sicherung guter wissenschaftlicher Praxis in der jeweils gültigen Fassung beachtet habe.

Karlsruhe, den 22.12.2017

Abstract

Most of the multilevel hypergraph partitioner implements the *FM* heuristic [16] to improve a given k -way partition. In each step, the *FM* algorithm moves a node with maximum *gain* from its current block to an other block. The *gain* of a node is the difference of an objective function before and after moving the node. The algorithm is easy to implement, flexible in adapting to different optimization functions and relatively fast [49]. However, a move is performed locally and greedily. Therefore, *FM* tends to find local optimal solutions which quality heavily depends on the *initial partition* [14]. Consequently, the probability to find a good approximation of a global optimum drops significantly if we partition large hypergraphs. Moreover, if we have a large number of moves with equal gain, the quality of the solution highly depends on random choices of the algorithm which makes the gain *meaningless* [28, 33]. Especially in large hyperedges, we often have many moves of vertices which gain is equal to zero if e.g., *cut* is the objective function.

Techniques derived from the *max-flow min-cut* theorem [17] to improve a k -way partition of a (hyper)graph can be seen as counterpart to the *FM* heuristic. They act *globally* and are not *move-based*. In this thesis, we present an alternative *local search* technique based on *Max-Flow-Min-Cut* computations for balanced direct k -way hypergraph partitioning. We integrate the work into the multilevel hypergraph partitioner *KaHyPar* [22]. The framework is inspired by the work of Sanders and Schulz [42] who successfully showed that *flow-based local search* in combination with the *FM* algorithm significantly improve the quality of partitions in a multilevel graph partitioner. We present several techniques to sparsify the state-of-the-art hypergraph flow network representation and show that maximum flow algorithms are up to a factor of 3 faster on our network in practice. Further, we show how to improve the *connectivity* metric of a k -way partition by solving a flow problem only on a subset of the vertices. We tested our new approach on a large benchmark set with 3222 instances. In comparison to 5 different systems, our new configuration outperforms the tested state-of-the-art hypergraph partitioner on 70% of the instances. In comparison to the latest version of *KaHyPar* our new approach produces 2% better quality by a performance slowdown only by a factor of 2. Moreover, our partitioner has a comparable running time to the direct k -way partitioner of *hMetis* and outperforms it on 82% of the benchmark instances with a 14.71% better solution quality.

Acknowledgements

A good mentor is important for staying motivated and to have the feeling to work on something meaningful. Since my bachelor thesis, I work together with Sebastian Schlag in the research area of *Hypergraph Partitioning*. The decision to further work on this topic was not only a matter of interest, it was mainly a decision based on our outstanding interpersonal working relationship. Our intellectual discussions were characterized by the right balance of fun and the necessary seriousness to work towards a common goal. I think here is the right place to thank you for the endless time, which I have spent in your office over the last three years and which have made me a better computer scientist.

Further, I would like to thank my girl friend Alessa Dreixler. To be together with a computer scientist could be sometimes very complicated and exhausting. Especially, if anger and frustration dominated my working day. With her understanding and motivation, I could continue every morning with the same energy and passion as the day before.

Contents

1. Introduction	7
1.1. Problem Statement	8
1.2. Contributions	8
1.3. Outline	9
2. Preliminaries	10
2.1. Graphs	10
2.2. Flows and Applications	10
2.3. Hypergraphs	12
2.4. Hypergraph Partitioning	13
3. Related Work	15
3.1. Maximum Flow Algorithms	15
3.1.1. Augmenting-Path Algorithms	15
3.1.2. Push-Relabel Algorithm	15
3.2. Modeling Flows on Hypergraphs	17
3.3. Flow-based Local Search on Graphs	18
3.3.1. Balanced Bipartitioning	18
3.3.2. Adaptive Flow Iterations	19
3.3.3. Most Balanced Minimum Cut	20
3.3.4. Active Block Scheduling	20
3.4. Hypergraph Partitioning	21
3.4.1. Multilevel Paradigm	21
3.4.2. n -Level Hypergraph Partitioning	22
4. Hypergraph Flow Networks	24
4.1. Removing Hypernodes via Clique-Expansion	24
4.2. Low-Degree Hypernodes	27
4.3. Removing Graph Hyperedges	27
4.4. Combining Techniques	29
5. Max-Flow-Min-Cut Refinement Framework	32
5.1. Source and Sink Configuration	32
5.2. Most Balanced Minimum Cuts on Hypergraphs	38
5.3. A direct k -way Flow-Based Refinement Framework	39
6. Experimental Results	42
6.1. Instances	42
6.2. System and Methodology	42
6.3. Flow Algorithms and Networks	43
6.4. Setup of the direct k -way Flow-Based Refinement	44
6.5. Speed-Up Heuristics	47
6.6. Comparison with other Hypergraph Partitioner	47
7. Conclusion	50
7.1. Future Work	50
A. Benchmark Instances	56
A.1. Parameter Tuning Benchmark Set	56

CONTENTS

A.2. Benchmark Subset	56
A.3. Full Benchmark Set	57
A.4. Excluded Test Instances	57
B. Detailed Flow Network and Algorithm Evaluation	59
C. Effectiveness Tests for Flow Configurations	60
D. Detailed Speedup Heuristic Evaluation	61
E. Detailed Comparison with other Systems	62

1. Introduction

Hypergraphs are a generalization of graphs, where each (hyper)edge can connect more than two (hyper)nodes. The k -way hypergraph partitioning problem is to partition the vertices of a hypergraph into k disjoint, non-empty blocks such that the size of each block is smaller than $1 + \epsilon$ times the average block size, while the goal is to simultaneously minimize an objective function.

Classical application areas can be found in *VLSI* design, parallelization of the Sparse Matrix-Vector Product and simplifying *SAT* formulas [26, 33, 37]. The goal in *VLSI* design is to partition a circuit into smaller units such that the wires between the gates are as short as possible [8]. A wire can connect more than two gates, therefore a hypergraph models a circuit more accurately than a graph. In *SAT* solving, hypergraph partitioning is used to decompose a formula into smaller subformulas, which can be solved more easily [33].

Hypergraph partitioning is an NP-hard problem [31] and it is even NP-hard to find a good approximation [7]. The most common heuristic used in state-of-the-art hypergraph partitioners is the *multilevel paradigm* [9, 22, 26]. First, a sequence of smaller hypergraphs is generated by contracting a set of matchings between hypernode pairs or clusters in each step (*coarsening phase*). If the hypergraph is small enough, we can use expensive heuristics to *initial partition* it into k blocks. Afterwards, the sequence of smaller hypergraphs is *uncontracted* in reverse order and, at each level, a *local search* heuristic is used to improve the quality of the partition (*refinement phase*).

There exist several *local search* heuristics for improving hypergraph partitions. One algorithm used in the state-of-the-art multilevel hypergraph partitioners is the *Fiduccia-Mattheyses* heuristic (FM). The FM algorithm maintains gain values (according to a objective function) of moving a node from its current block to another block [16]. A move is performed, if its gain value is maximum among all possible moves. The FM heuristic is generally intuitive, flexible in adapting to different optimization objectives, easy to implement and relatively fast [49]. However, a move is performed without a *global* view on the problem instance and *greedily*. Consequently, the algorithm tends to find locally optimal solutions, which quality heavily depends on the initial partition [14]. Therefore, multiple runs are needed to find a solution close to the global optimum. The probability of finding a good approximation of the global optimum significantly drops if we partition large hypergraphs [14]. Executing FM in the multilevel context partially solve the disadvantage for large benchmarks. A move of an vertex in a *coarsened* hypergraph corresponds to a movement of a subset of the hypernodes in the *original* hypergraph which allows a more effective exploration of the solution space [37]. Further, a move of a node only influences the gain function if the state of an incident hyperedge changes *immediately* after a move. If a hyperedge contains vertices from two different blocks, where only one hypernode is contained in the first and all remainings are in the second block, then a move of that node contributes to the gain if the objective is e.g., *cut* (sum of the weights of hyperedges which contains vertices of more than one block). Especially for large hyperedges, we often have to *move* a sequence of nodes such that a single move of a node finally contributes to the gain. Therefore, the gain of most vertices is equal to zero in such cases [33]. Krishnamurthy [28] points out that the quality in such situations highly depends on random choices made within the algorithm. Therefore, he enhanced the FM algorithm with a look-ahead scheme such that in case of ties one can incorporate *future gains* into the decision [28]. However, the *forecast* is limited by a predefined parameter.

FM-based *local search* algorithms have the above-mentioned disadvantages, because they are *move-based* and only incorporate *local* informations about the structure of the problem. Finding a balanced global minimum cut of a (hyper)graph is NP-hard, but if we ask for a minimum

cut separating two vertices s and t the problem becomes solvable in polynomial time [15]. The well-known *max-flow min-cut* theorem [17] established an analogy between the maximum flow from a source s to a sink t and the minimum cut separating s and t in a graph. *Flow-based* approaches are not *move-based* and incorporate the *global* structure of the problem. Therefore, it overcomes the drawbacks of the *FM* algorithm. However, it was overlooked for a long time because it was perceived as computationally expensive and impractical for (hyper)graph partitioning [32].

Recently, several algorithms to obtain a balanced bipartition of a hypergraph were developed [32, 38, 48]. The algorithms can be seen as an alternative to the *multilevel paradigm*. The impact of a *flow-based local search* algorithm on the solution quality of a multilevel hypergraph partitioner has not been studied yet. Moreover, a balanced k -way hypergraph partition is only obtainable by applying the bipartitioning algorithm recursively [48]. Sanders and Schulz [42] successfully integrated a *flow-based refinement* algorithm in their multilevel *graph* partitioner. The algorithm is also applicable to the more complicated direct k -way partitioning case. In general, their basic approach is to extract a subgraph around the cut of a bipartition and configure the source and sink sets of the flow problem such that a maximum flow calculation on the subgraph leads to a smaller cut on the original graph. They combine the strength of *flow-based* and *FM local search* by executing both algorithms alternating throughout the multilevel hierarchy. As a result their multilevel graph partitioner produces the best partitions for a wide range of graph partitioning benchmarks.

1.1. Problem Statement

Motivated by the successfull integration of a *flow-based local search* algorithm in the multilevel graph partitioner *KaFFPa* of Sanders and Schulz [42] to obtain balanced k -way partitions, this thesis investigates the integration of such an approach into the multilevel hypergraph partitioner *KaHyPar* [22].

In the first step, we have to find an appropriate model of a hypergraph as flow network. Afterwards, we want to improve a given bipartition of a hypergraph with a *Max-Flow-Min-Cut* computation by using the flow network of the previous step. This work are the theoretical foundations for developing a *flow-based local search* algorithm which works in a multilevel hypergraph partitioner and improves a given balanced k -way partition. The last step is to integrate the framework into the n -level hypergraph partitioner *KaHyPar* [22] and evaluate the performance on a large benchmark set in comparison to different state-of-the-art multilevel hypergraph partitioners. A major goal of this work is to outperform the latest version of *KaHyPar* on most of the benchmark instances and simultaneously ensure that the running time is only within a constant factor slower.

1.2. Contributions

We present several techniques to sparsify the state-of-the-art hypergraph flow network modeling approach proposed by Lawler [30]. Our experiments indicate that maximum flow algorithms are up to a factor of 3 faster with our new network. The theoretical results, which leads to the presented sparsification techniques, are of independent interest and can also be applied on general flow networks. Further, we show how to configure a flow problem on a subhypergraph of H such that a maximum (S, T) -flow yields an improved cut of a given bipartition. We choose S and T in a way such that the value of the cut of H after a *Max-Flow-Min-Cut* computation on a subhypergraph can be calculated with the value of a maximum (S, T) -flow. Our *flow-based local search* framework is inspired by algorithmic ideas of Sanders and Schulz [42]. However, we

generalize many results of their work such that they are applicable on hypergraph partitioning. Further, we implement several heuristics which might prevent unpromising *Max-Flow-Min-Cut* computations throughout the multilevel hierarchy and show that they speed-up the framework by factor of 2 while maintaining the quality of the solutions on average. We integrate our *flow-based local search* algorithm into the n -level hypergraph partitioner *KaHyPar* and show that *flow-based refinement* in combination with the *FM* algorithm produces the best partitions on a majority of real world benchmarks in comparison to other state-of-the-art hypergraph partitioners. Compared to 5 different systems we achieve the best partitions on 70% of 3222 benchmark instances. In comparison to the latest version of *KaHyPar*, our new approach produces solutions that are 2% better on average, while only incurring a slowdown by a factor of 2. Moreover, our partitioner has a comparable running time to the direct k -way version of *hMetis* and outperforms it on 82% of the benchmark instances.

1.3. Outline

We first introduce necessary notations and summarize related work in Sections 2 and 3. Afterwards, we describe techniques to sparsify the flow network proposed by Lawler [30] in Section 4. In Section 5 we present our source and sink set modeling approach and describe the integration of our *flow-based refinement* framework into the n -level hypergraph partitioner *KaHyPar*. The experimental evaluation of our algorithm is presented in Section 6. Section 7 concludes this thesis.

2. Preliminaries

2.1. Graphs

Definition 2.1. A directed weighted graph $G = (V, E, c, \omega)$ is a set of nodes V and a set of edges E with a node weight function $c : V \rightarrow \mathbb{R}_{\geq 0}$ and an edge weight function $\omega : E \rightarrow \mathbb{R}_{\geq 0}$. An edge $e = (u, v)$ is a relation between two nodes $u, v \in V$.

Two vertices u and v are *adjacent*, if there exists an edge $(u, v) \in E$. Two edges e_1 and e_2 are *incident* to each other if they share a node. $N(v)$ denotes the set of all *adjacent* nodes of v . The *degree* of a node v is $d(v) = |N(v)|$.

Definition 2.2. Given a directed graph $G = (V, E)$. A contraction of two nodes u and v results in a new graph $G_{(u,v)} = (V \setminus \{v\}, E')$, where each edge of the form (v, w) or (w, v) in E is replaced with an edge (u, w) or (w, u) in E' .

A *path* $P = (v_1, \dots, v_k)$ is a sequence of nodes, where for each $i \in [1, k - 1] : (v_i, v_{i+1}) \in E$. A *cycle* is a path $P = (v_1, \dots, v_k)$ with $v_1 = v_k$. A *strongly connected component* $C \subseteq V$ is a set of nodes where for each $u, v \in C$ exists a path from u to v . We can enumerate all *strongly connected components* (*SCC*) in a directed graph G with a linear time algorithm proposed by Tarjan [45]. A directed graph G without any *cycles* is called *directed acyclic graph* (*DAG*). On such graphs we can define a *topological order* $\gamma : V \rightarrow \mathbb{N}_+$ such that for each $(u, v) \in E : \gamma(u) < \gamma(v)$. A *topological order* of a *DAG* can be found in linear time with Kahn's algorithm [25]. We can transform a general directed graph G into a *DAG* if we contract each *strongly connected component*. All concepts are illustrated in Figure 1.

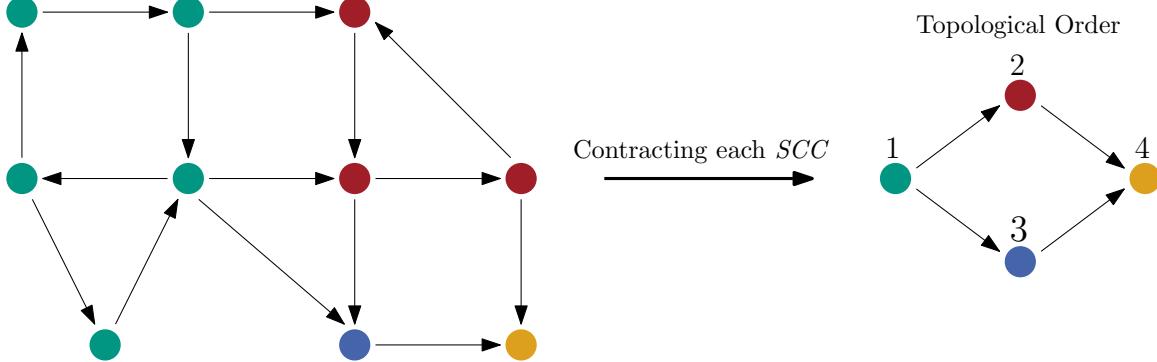


Figure 1: Example of *Strongly Connected Components* in a directed graph and a *Topological Order* of a *Directed Acyclic Graph*. Each *SCC* is marked with the same color.

Definition 2.3. Let $G_{V'} = (V', E_{V'}, c, \omega)$ be the subgraph of a graph G induced by $V' \subseteq V$ with $E_{V'} = \{(u, v) \in E \mid u, v \in V'\}$.

2.2. Flows and Applications

Given a graph $G = (V, E, u)$ with capacity function $u : E \rightarrow \mathbb{R}_+$ and a source $s \in V$ and a sink $t \in V$, the maximum flow problem is to find the maximum amount of flow from s to t in G . A flow is a function $f : E \rightarrow \mathbb{R}_+$, which have to satisfy the following constraints:

- (i) $\forall (u, v) \in E : f(u, v) \leq u(u, v)$ (capacity constraint)
- (ii) $\forall v \in V \setminus \{s, t\} : \sum_{(u,v) \in E} f(u, v) = \sum_{(v,u) \in E} f(v, u)$ (conservation of flow constraint)

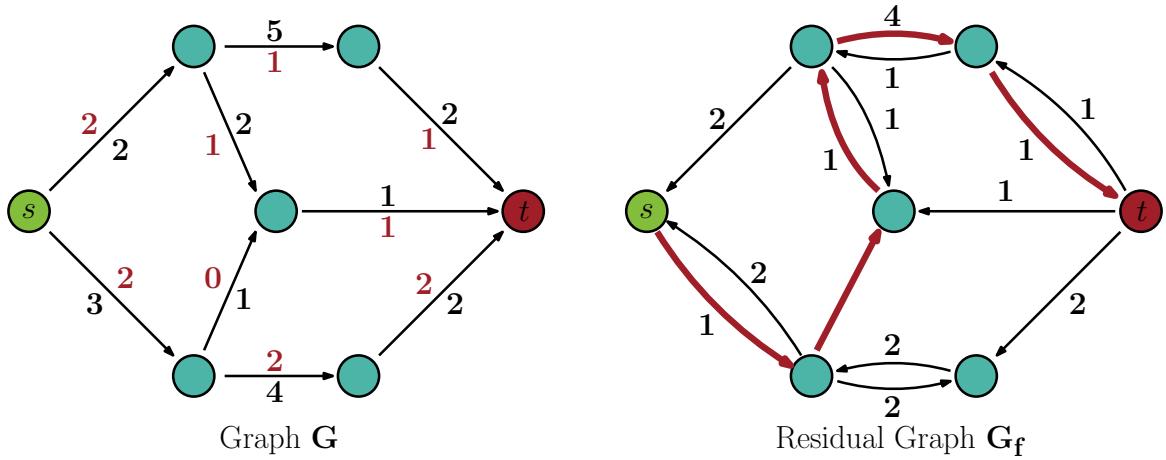


Figure 2: Illustrates concepts related to the maximum flow problem. A flow function f (red values) from s to t of a graph G is shown on the left side. The corresponding *residual graph* G_f with its *residual capacities* (black values) is illustrated on the right side. The red highlighted path is an *augmenting path*.

The capacity constraint restricts the flow on edge (u, v) by its capacity $u(u, v)$. Whereas the conservation of flow constraint ensures that the amount of flow entering a node $v \in V \setminus \{s, t\}$ is the same as leaving a node. The value of the flow is defined as $|f| = \sum_{(s,v) \in E} f(s, v) = \sum_{(v,t) \in E} f(v, t)$. A flow f is maximal, if there exists no other flow f' with $|f'| > |f|$.

Further, we define the *residual graph* G_f and the *residual capacity* r_f of a flow function f on graph G . The *residual capacity* $r_f : V \times V \rightarrow \mathbb{R}_+$ is defined as follows:

- (i) $\forall (u, v) \in E : r_f(u, v) = u(u, v) - f(u, v)$
- (ii) $\forall (u, v) \in E : \text{If } f(u, v) > 0 \text{ and } u(v, u) = 0, \text{ then } r_f(v, u) = f(u, v)$

For an edge $e = (u, v) \in E$ the residual capacity $r_f(u, v)$ is the remaining amount of flow which can be send over edge e . For each reverse edge $\overleftarrow{e} \notin E$ the residual capacity $r_f(\overleftarrow{e})$ is the amount of flow which is send over e . The *residual graph* $G_f = (V, E_f, r_f)$ is the network containing all $(u, v) \in V \times V$ with $r_f(u, v) > 0$. More formally $E_f = \{(u, v) \in V \times V \mid r_f(u, v) > 0\}$. An *augmenting path* $P = \{v_1, \dots, v_k\}$ is a path in G_f with $v_1 = s$ and $v_k = t$ [15]. Figure 2 illustrates all presented concepts.

The *Max-Flow-Min-Cut-Theorem* is fundamental for many applications related to the maximum flow problem [17].

Theorem 2.1. *The value of a maximum (s, t) -flow obtainable in a graph G is equal with the minimum-weight cutset in G separating s and t .*

Let f be a maximum (s, t) -flow in a graph $G = (V, E, \omega)$ with $s \in V$ and $t \in V$. Further, let A be the set containing all $v \in V$, which are *reachable* from s in G_f . A node v is *reachable* from a node u if there exists a path from u to v . Then the set of all cut edges between the bipartition $(A, V \setminus A)$ is a minimum-weight (s, t) -cutset [18]. A can be calculated with a *BFS* in G_f starting from s .

We can solve with maximum flows many related problems like e.g., maximum bipartite-matching, number of edge- or vertex-disjoint paths in a graph or to find a minimum-weight vertex separator. Solutions for those problems sometimes involve a transformation T of the graph G into a flow network $T(G)$, such that the *Max-Flow-Min-Cut-Theorem* is applicable. A problem essential for this work is to find a minimum-weight (s, t) -vertex separator in a graph $G = (V, E, c)$ with $c : V \rightarrow \mathbb{R}_+$.

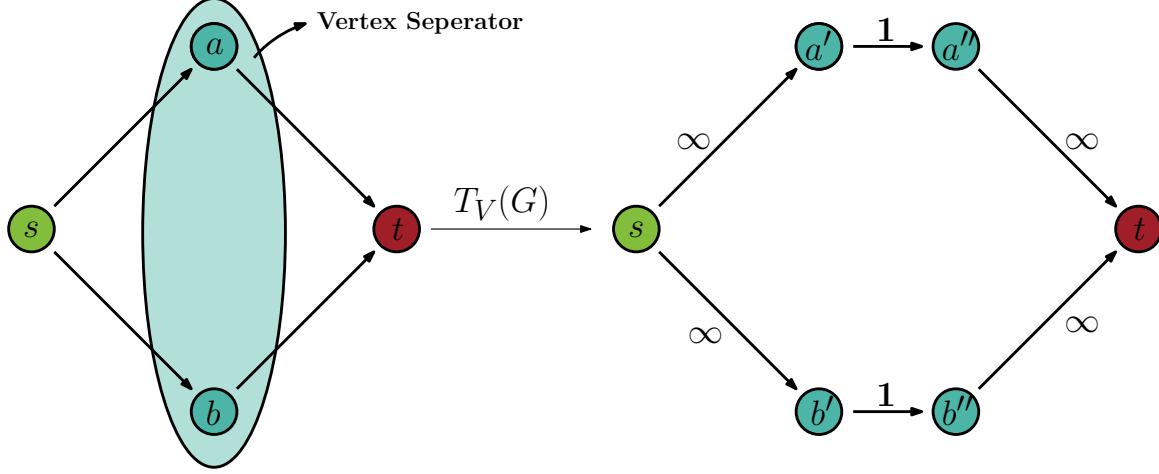


Figure 3: Illustration of the vertex separator problem and the flow network $T_V(G)$ in which we can find a minimum vertex separator.

Definition 2.4. Let $G = (V, E, c)$ be a graph with $c : V \rightarrow \mathbb{R}_+$. $S \subseteq V$ is a vertex separator for non-adjacent vertices $s \in V$ and $t \in V$ if the removal of S from graph G separates s and t (s not reachable from t). A vertex separator S is a minimum-weight (s, t) -vertex separator, if for all (s, t) -vertex separators $S' \subseteq V$ follows that $c(S) \leq c(S')$.

We can calculate a minimum-weight (s, t) -vertex separator with a maximum flow calculation on the following flow network [47]:

Definition 2.5. Let T_V be a transformation of a graph $G = (V, E, c)$ into a flow network $T_V(G) = (V_V, E_V, u_V)$ (with $u_V : E_V \rightarrow \mathbb{R}_+$). T_V is defined as follows:

- (i) $V_V = \bigcup_{v \in V} \{v', v''\}$
- (ii) $\forall v \in V$ add a directed edge (v', v'') with capacity $u_V(v', v'') = c(v)$
- (iii) $\forall (u, v) \in E$ add two directed edges (u'', v') and (v'', u') with capacity $u_V(u'', v') = u_V(v'', u') = \infty$.

The vertex separator problem and transformation $T_V(G)$ is illustrated in Figure 3. Obviously, no edge between two adjacent nodes in G can be in a minimum-capacity (s, t) -cutset of $T_V(G)$, because for all those edges the capacity is ∞ . Therefore, the cutset must consist of edges of the form (v', v'') . A minimum-weight (s, t) -vertex separator can be calculated by finding a maximum (s, t) -flow of $T_V(G)$ and the corresponding minimum (s, t) -cutset [34].

Given a set of sources S and sinks T . The *multi-source multi-sink* maximum flow problem is to find a maximum flow f from all source nodes $s \in S$ to all sink nodes $t \in T$. We can transform such a problem into a *single-source single-sink* problem by adding two additional nodes s and t . We add a directed edge from s to all source nodes $s' \in S$ and for all sink nodes $t' \in T$ a directed edge to t with capacity $u(s, s') = u(t', t) = \infty$.

2.3. Hypergraphs

Definition 2.6. An undirected weighted hypergraph $H = (V, E, c, \omega)$ is a set of hypernodes V and a set of hyperedges E with a hypernode weight function $c : V \rightarrow \mathbb{R}_{\geq 0}$ and a hyperedge weight function $\omega : E \rightarrow \mathbb{R}_{\geq 0}$. A hyperedge e is a subset of V (formally: $\forall e \in E : e \subseteq V$).

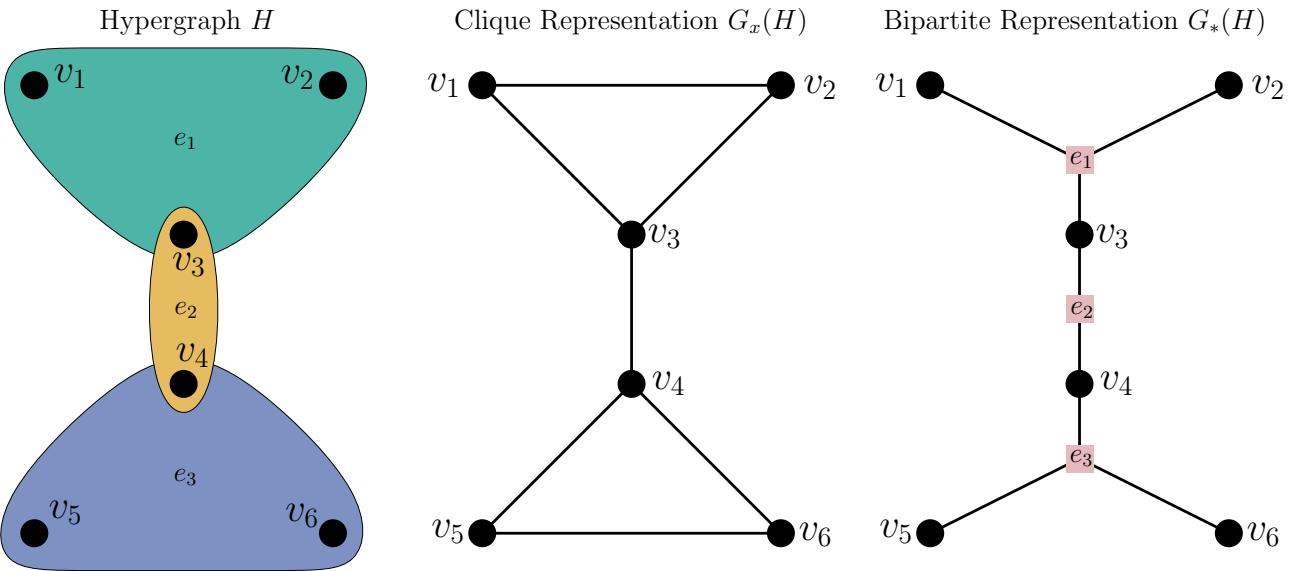


Figure 4: Example of a hypergraph H and its two corresponding graph representations.

A hypergraph generalizes a graph by extending the definition of an edge, which can contain more than two nodes. Hyperedges are also called *nets* and hypernodes are also called *vertices*. A vertex contained in a net is called *pin*. For a subset $V' \subseteq V$ and $E' \subseteq E$ we define

$$c(V') = \sum_{v \in V'} c(v)$$

$$\omega(E') = \sum_{e \in E'} \omega(e)$$

A vertex v is *incident* to a hyperedge e if $v \in e$. Two vertices u and v are *adjacent*, if there exists an $e \in E$ such that $u \in e$ and $v \in e$. $I(v)$ denotes the set of all *incident* nets of v . The *degree* of a hypernode v is $d(v) = |I(v)|$. The size of a net e is its cardinality $|e|$.

Definition 2.7. Let $H_{V'} = (V', E_{V'}, c, \omega)$ be the subhypergraph of a hypergraph H induced by $V' \subseteq V$ with $E_{V'} = \{e \cap V' \mid e \in E : e \cap V' \neq \emptyset\}$.

A hypergraph $H = (V, E, c, \omega)$ can be represented as an undirected graph. There are two standard transformations, called *clique* and *bipartite* representation [24]. The *clique* graph $G_x(H) = (V, E_x)$ models each net e as a clique between its pins. The *bipartite* graph $G_*(H) = (V \cup E, E_*)$ contains all hypernodes and hyperedges as nodes and connects each net e with an undirected edge $\{e, v\}$ to all its pins $v \in e$. The two transformations are illustrated in Figure 4.

2.4. Hypergraph Partitioning

Definition 2.8. A k -way partition of a hypergraph H is a partition of its hypernodes into k disjoint blocks $\Pi = \{V_1, \dots, V_k\}$ such that $\bigcup_{i=1}^k V_i = V$ and $V_i \neq \emptyset$.

For a k -way partition $\Pi = \{V_1, \dots, V_k\}$, we define the *connectivity set* of a hyperedge e with $\Lambda(e, \Pi) = \{V_i \in \Pi \mid V_i \cap e \neq \emptyset\}$. The *connectivity* of a net e is $\lambda(e, \Pi) = |\Lambda(e, \Pi)|$. A hyperedge e is *cut*, if $\lambda(e, \Pi) > 1$. $E(\Pi) = \{e \mid \lambda(e, \Pi) > 1\}$ is the set of all *cut* nets. We say two blocks V_i and V_j are adjacent, if there exists a hyperedge e with $V_i \in \Lambda(e, \Pi)$ and $V_j \in \Lambda(e, \Pi)$. A k -way partition is ϵ -balanced if each block $V_i \in \Pi$ satisfies the *balance constraint* $c(V_i) \leq (1 + \epsilon) \lceil \frac{c(V)}{k} \rceil$.

Definition 2.9. *The k -way hypergraph partitioning problem is to find an ϵ -balanced k -way partition Π of a hypergraph H such that a certain objective function is minimized.*

There exists several objective functions in the hypergraph partitioning context. The most popular objective function is the cut metric (especially for *graph partitioning*), which is defined as

$$\omega_H(\Pi) = \sum_{e \in E(\Pi)} \omega(e)$$

The goal is to minimize the weight of all *cut* hyperedges. Another important metric for this work is the $(\lambda - 1)$ -metric or *connectivity* metric, which is defined as

$$(\lambda - 1)_H(\Pi) = \sum_{e \in E} (\lambda(e) - 1)\omega(e)$$

The idea behind this function is to minimize the *connectivity* of all hyperedges.

Definition 2.10. *We define for a k -way partition $\Pi = \{V_1, \dots, V_k\}$ of a hypergraph H the quotient graph $Q = (\Pi, E')$ which contains an edge between each pair of adjacent blocks of Π . More formally, $E' = \{(V_i, V_j) \mid \exists e \in E : V_i, V_j \in \Lambda(e, \Pi)\}$*

3. Related Work

3.1. Maximum Flow Algorithms

In Section 2.2 we introduce the concept of flows in a network. We will now present two algorithms to solve the maximum flow problem.

3.1.1. Augmenting-Path Algorithms

An *augmenting path* $P = \{v_1, \dots, v_k\}$ is a path in G_f with $v_1 = s$ and $v_k = t$ [15]. Figure 5 illustrates such a path. Since all $(v_i, v_{i+1}) \in G_f$ it follows that $r_f(v_i, v_{i+1}) > 0$. Therefore, we can increase the flow on all edges (v_i, v_{i+1}) by $\Delta f = \min_{i \in [1, \dots, k-1]} r_f(v_i, v_{i+1})$. It can be shown that f is not a maximum flow if an *augmenting path* exists in G_f [15].

One way to calculate a maximum flow f is to find *augmenting paths* in G_f as long as there exists one. The algorithm was established by Ford and Fulkerson [17] and consists of two phases. First, we search for an *augmenting path* $P = \{v_1, \dots, v_k\}$ from s to t , e.g., with a simple *DFS*. Afterwards, we increase the flow on each edge (v_i, v_{i+1}) by Δf and decrease the flow on each reverse edge (v_{i+1}, v_i) by Δf . If the capacities are integral, the algorithm always terminates. Since we can find an *augmenting path* in G_f with a simple *DFS* in $\mathcal{O}(|V| + |E|)$ and increase the flow on every path by at least one, the running time of the algorithm can be bounded by $\mathcal{O}(|E||f_{max}|)$. We can construct instances, where the running time is $\mathcal{O}(|E||f_{max}|)$ or even the maximum flow $|f_{max}|$ is exponential in the problem size [15].

Edmond and Karp [15] improved Ford & Fulkerson's algorithm by increasing the flow along an *augmenting path* of minimal length. The shortest path from s to t in a graph with unit lengths can be found with a simple *BFS*. It can be shown, that the total number of *augmentations* is $\mathcal{O}(|V||E|)$. The running time of Edmond & Karp's maximum flow algorithm is $\mathcal{O}(|V||E|^2)$. A sample execution of the algorithm is presented in Figure 5.

3.1.2. Push-Relabel Algorithm

Goldberg and Tarjan [20] implemented a maximum flow algorithm not based on finding an *augmenting path* in the *residual graph*. The idea is to maintain a *preflow* during the execution of the algorithm which satisfies the capacity constraints, but only a weakened form of the conservation of flow constraint:

$$\forall v \in V \setminus \{s, t\} : \sum_{u \in V} f(v, u) \leq \sum_{u \in V} f(u, v)$$

The algorithm maintains a *distance labeling* $d : V \rightarrow \mathbb{N}$ and an *excess function* $e_f : V \rightarrow \mathbb{N}$. The *distance labeling* satisfies the following conditions: $d(s) = |V|$, $d(t) = 0$ and for each $(u, v) \in E_f$, $d(u) \leq d(v) + 1$. We say an residual edge (u, v) is *admissible* if $d(u) = d(v) + 1$. A node v is *active* if $v \notin \{s, t\}$ and $e_f(v) > 0$.

Initially, all *labels* and *excess* values are set to zero except source node s will be set to $d(s) = 1$ and $e_f(s) = \infty$. For each *active* node u the algorithm performs two update operations, called *push* and *relabel*. The first operation pushes flow over each *admissible* edge (u, v) . After a *push* $e_f(u) = e_f(u) - \min(e_f(u), r_f(u, v))$ and $e_f(v) = e_f(v) + \min(e_f(u), r_f(u, v))$. If there is no *admissible* edge, a *relabel* operation is performed, which replaces $d(u)$ by $\min_{(u,v) \in E_f} d(v) + 1$. The algorithm terminates, if none of the nodes is *active*. The worst case complexity of the algorithm is $\mathcal{O}(n^3)$. The running time can be reduced to $\mathcal{O}(n^2 \log n)$ with *Dynamic Trees* [20, 44], but this implementation is not practical due to a large hidden constant factor.

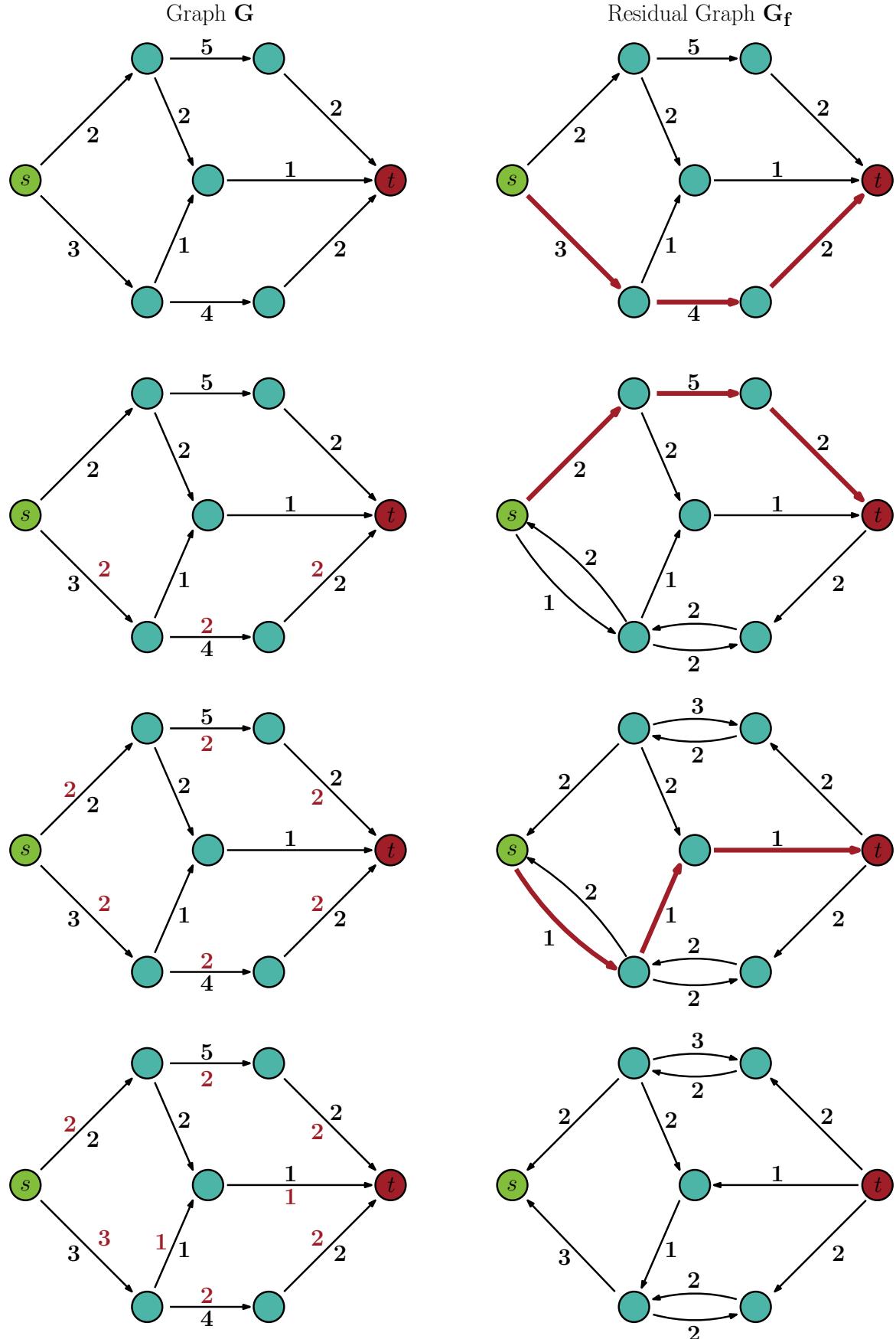


Figure 5: Execution of Edmond & Karps maximum flow algorithm [15]. The network G with its capacities c (black values) and flow f (red values) is illustrated on the left side. The residual graph G_f with its *residual capacities* r_f (black values) is presented on the right side. In each step the current *augmenting path* in G_f is highlighted by a red path.

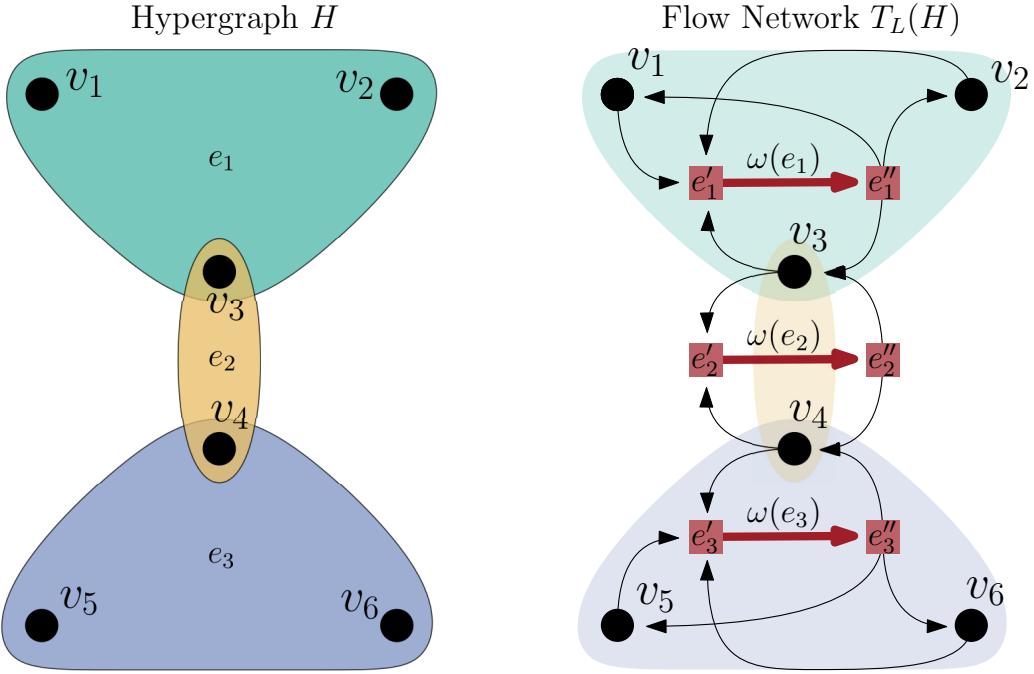


Figure 6: Transformation of a hypergraph into an equivalent flow network by Lawler [30]. Note, capacity of the black edges in the flow network is ∞ .

The *push-relabel* algorithm is one of the fastest maximum flow algorithms in practice because there exist several speed-up techniques. The first one is the *global relabeling* heuristic which frequently updates the *distance labels* by computing the shortest path in the residual graph from all nodes to the sink [11]. This can be done with a backward *BFS* in linear time. This technique is performed periodically, e.g., after every n relabeling.

The second heuristic is the *gap heuristic* [10, 13]. If at a particular stage of the algorithm there is no node u with $d(u) = g < n$, then for each node v with $g < d(v) < n$ the sink is not reachable anymore. Therefore, we can increase the *distance label* of all those nodes to n . To implement this heuristic, we maintain a linked list of nodes with distance label i .

3.2. Modeling Flows on Hypergraphs

Consider the *bipartite graph* representation $G_*(H)$ of a hypergraph H (see Section 2.3). Hu and Moerder [24] introduce node capacities in $G_*(H)$. Each hyperedge e has a capacity equal to $\omega(e)$ and each hypernode has infinite capacity. Further, they show that a minimum-weight (s, t) -vertex separator in $G_*(H)$ is equal with a minimum-weight (s, t) -cutset of a hypergraph H . Finding such a separator is a flow problem and can be calculated with the flow network $T_L(H)$ presented by Lawler [30]:

Definition 3.1. Let T_L be the transformation of a hypergraph $H = (V, E, c, \omega)$ into a flow network $T_L(H) = (V_L, E_L, u_L)$ proposed by Lawler [30]. $T_L(H)$ is defined as follows:

- (i) $V_L = V \cup \bigcup_{e \in E} \{e', e''\}$
- (ii) $\forall e \in E$ we add a directed edge (e', e'') with capacity $u_L(e', e'') = \omega(e)$
- (iii) $\forall v \in V$ and $\forall e \in I(v)$ we add two directed edges (v, e') and (e'', v) with capacity $u_L(v, e') = u_L(e'', v) = \infty$.

An example of this transformation is shown in Figure 6. $T_L(H)$ is nearly equivalent to the transformation $T_V(G)$ described in Definition 2.5 except that we do not have to split the hypernodes

$v \in V$. A hypernode cannot be in a minimum-capacity (s, t) -vertex separator because each $v \in V$ has infinity capacity [24]. Therefore, a minimum-capacity (s, t) -cutset of $T_L(H)$ is equal to a minimum (s, t) -vertex separator of $G_*(H)$. The resulting graph $T_L(H)$ has $|V_L| = 2|V| + |E|$ nodes and $|E_L| = 2(\bar{e} + 1)|E|$ edges, where \bar{e} is the average size of a hyperedge [40]. Using *Edmond-Karps* maximum flow algorithm (see Section 3.1.1) on flow network $T_L(H)$ takes time $\mathcal{O}(|V|^2|E|^2)$ [30].

A minimum-weight (s, t) -cutset of H can be found by simply mapping the minimum-capacity (s, t) -cutset to their corresponding hyperedges in H (see Section 2.2). The minimum-weight (s, t) -bipartition are all vertices $v \in V$ *reachable* from s in the *residual graph* of $T_L(H)$ and the counterpart are all hypernodes not *reachable* from s .

In Figure 7 we illustrate the structure of $T_L(H)$ and demonstrate what happens after we *augment* along a path in the Lawler-Network. This figure can be used as a reference if you need an illustration of techniques used in the proofs of Section 4.

In this thesis, we often have to mix up nodes and edges of H and $T_L(H)$. If we use $v \in V_L$, there also exists a corresponding $v \in V$. v can be used in both contexts. For all $e \in E$ there exists two corresponding nodes $e', e'' \in V_L$. e' is called *incoming hyperedge node* and e'' is called *outgoing hyperedge node*. In some cases we need to treat $e', e'' \in V_L$ the same way as their corresponding hyperedge $e \in E \Rightarrow e'_1 \cap e''_2$ or $e''_1 \cap e'_2$ should be the same as $e_1 \cap e_2$.

3.3. Flow-based Local Search on Graphs

It seems natural to utilize maximum flow computations to improve the cut metric of a given partition of a graph. Lang and Rao [29] use an approach, called *Max-Flow Quotient-cut Improvement* (MQI), to improve the quality of a graph when metrics such as *expansion* or *conductance* are used. For a given bipartition (S, \bar{S}) , they find the best improvement among all bipartitions (S', \bar{S}') such that $S' \subset S$ by solving a flow problem. Andersen and Lang [4] suggested a flow-based improvement algorithm, called *Improve*, which works similar as MQI, but do not restrict the output of the partition to $S' \subset S$. However, both techniques can not guarantee that the resulting bipartition is balanced and only are applicable for $k = 2$.

Schulz and Sanders [42] integrate flow-based refinement algorithm in their *multilevel graph partitioner KaFFPa*. In general, they build a flow problem a region B around the cut and connect the *border* of B with the source resp. sink. B is defined in such a way that the flow computation yields a feasible cut according to the *balanced constraint*. Many ideas of this work are used in this thesis and adapted to hypergraphs. Therefore, we will give a detailed description of the concepts and advanced techniques to improve graph partitions.

3.3.1. Balanced Bipartitioning

Let (V_1, V_2) be a balanced bipartition of a graph $G = (V, E, c, \omega)$. Further, $P(v) = 1$, if $v \in V_1$ and $P(v) = 2$, otherwise. We will now explain how a given bipartition can be improved with flow computations. This technique can also be applied on a k -way partition by applying the approach on two adjacent blocks [42].

Let $\delta := \{u \mid \exists(u, v) \in E : P(u) \neq P(v)\}$ be the set of nodes around the cut of G . For a set $B \subseteq V$ we define its border $\delta B := \{u \in B \mid \exists(u, v) \in E : v \notin B\}$. The basic idea is to build a region B around all cut nodes δ of G and connect all nodes in $\delta B \cap V_1$ to the source node s and all nodes in $\delta B \cap V_2$ to the sink node t .

We can construct $B := B_1 \cup B_2$ with two *Breadth First Searches (BFS)*. One is initialized with all nodes $\delta \cap V_1$ and stops if $c(B_1)$ would exceed $(1 + \epsilon) \frac{c(V)}{2} - c(V_2)$. The second is initialized with all nodes $\delta \cap V_2$ and stops if $c(B_2)$ would exceed $(1 + \epsilon) \frac{c(V)}{2} - c(V_1)$. The two *BFSs* only

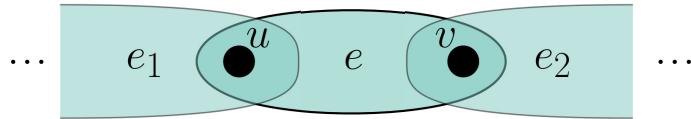
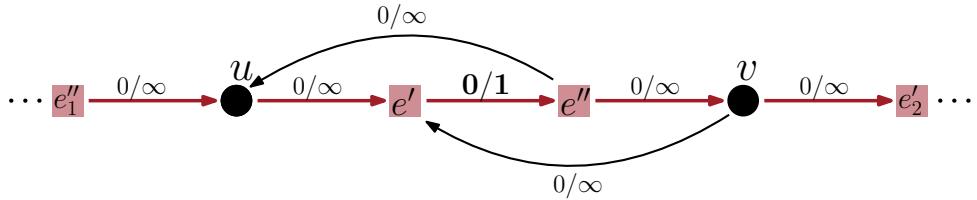
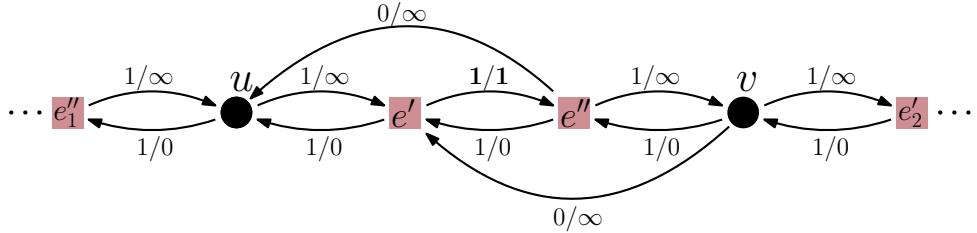
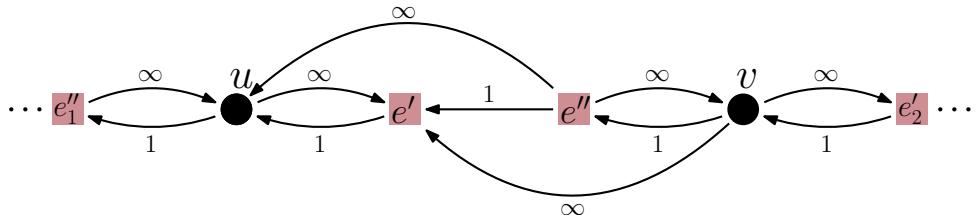
Hypergraph H Flow Network $T_L(H)$ Flow Network $T_L(H)$ after augmentingResidual Graph of $T_L(H)$ after augmenting

Figure 7: Illustration of the structure of a hyperedge e in $T_L(H)$ and the effect of augmenting along a path in this network. The labeling on an edge represents the flow $f(x,y)$ and the capacity $u(x,y)$ denoted with $f(x,y)/u(x,y)$. The labeling of the edges in the residual graph denotes the residual capacity $r_f(x,y)$. The red highlighted path represents an augmenting path.

touches nodes of V_1 resp. $V_2 \Rightarrow B_1 \subseteq V_1$ and $B_2 \subseteq V_2$. The constraints for the weights of B_1 and B_2 guarantees that the bipartition is still balanced after a *Max-Flow-Min-Cut* computation. Connecting s resp. t to all border nodes $\delta B \cap V_1$ resp. $\delta B \cap V_2$ ensures that a non-cut edge not contained in G_B is not a cut edge after assigning the minimum (s,t) -bipartition of subgraph G_B to G . This also yields the conclusion that each minimum (s,t) -cutset in G_B leads to a cut smaller or equal to the old cut of G . All concepts are illustrated in Figure 8.

3.3.2. Adaptive Flow Iterations

Sanders and Schulz [42] introduce several techniques to improve their basic approach. If the *Max-Flow-Min-Cut* computation on G_B leads to an improved cut, we can apply the method

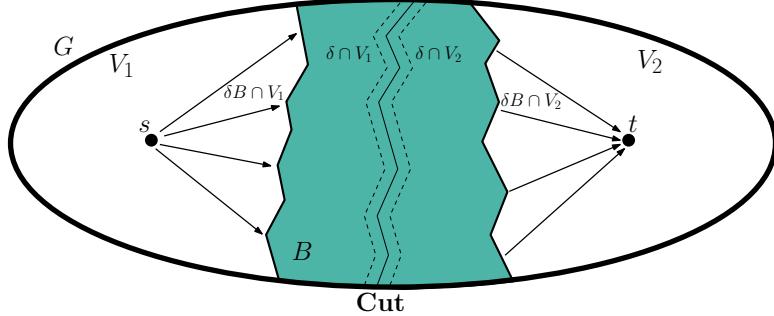


Figure 8: Configuration of a flow problem around the cut of graph G [4].

described in Section 3.3.1 again. An extension of this approach is to iteratively adapt the size of the flow problem based on the result of the maximum flow computation. We define $\epsilon' := \alpha\epsilon$ for a $\alpha \geq 1$ and let the size of B depend on ϵ' rather than on ϵ . If we find an improvement on G , we increase α to $\min\{2\alpha, \alpha'\}$ where α' is a predefined upper bound for α . If not, we decrease the size of α to $\max\{\frac{\alpha}{2}, 1\}$. This approach is called *adaptive flow iterations* [42].

3.3.3. Most Balanced Minimum Cut

Picard and Queyranne [39] show that all minimum (s, t) -cutsets are computable with one maximum (s, t) -flow computation. To understand the main theorem and the algorithm to compute all minimum (s, t) -cutsets we need the definition of a *closed node set* $C \subseteq V$ of a graph G .

Definition 3.2. Let $G = (V, E)$ be a graph and $C \subseteq V$. C is called a closed node set iff the condition $u \in C$ implies that for all edges $(u, v) \in E$ also $v \in C$.

A *closed node set* is illustrated in Figure 9. A simple observation is that all nodes on a cycle have to be in the same *closed node set* per definition. Therefore we can contract all *Strongly Connected Components* (SCC) of G with a linear time algorithm proposed by Tarjan [45] and sweep in reverse topological order over the contracted graph to enumerate all *closed node sets*. Note, if we contract all SCCs of G the resulting graph is a *Directed Acyclic Graph* (DAC). Therefore, a topological order exists.

With the Theorem of Picard and Queyranne [39] we can enumerate all minimum (s, t) -cuts of G with one maximum flow computation.

Theorem 3.1. There is a 1-1 correspondence between the minimum (s, t) -cuts of a graph and the closed node sets containing s in the residual graph of a maximum (s, t) -flow.

All *closed node sets* in the residual graph of G induce a minimum (s, t) -cutset on G . They can be calculated with the algorithm described above having the residual graph of G as input. The running time of the algorithm is $\mathcal{O}(|V| + |E|)$.

A common problem of the *adaptive flow iteration* approach (see Section 3.3.2) is that using a large α often leads to cuts in G that violate the balanced constraint. We can enumerate all minimum (s, t) -cutsets with one maximum flow computation and therefore have a higher probability to find a feasible partition after a *Max-Flow-Min-Cut* computation. We refer to this method as *Most Balanced Minimum Cut*.

3.3.4. Active Block Scheduling

Active Block Scheduling is a *quotient graph style refinement* technique for k -way partitions [23, 42]. The algorithm is organized in rounds and executes a two-way local improvement

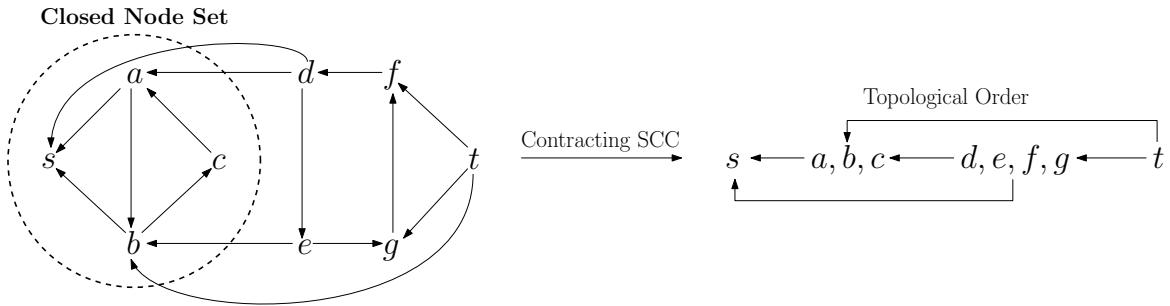


Figure 9: $C = \{s, a, b, c\}$ is a *closed node set* of graph G (left side). After contracting all *Strongly Connected Components*, we can enumerate all *closed node sets* of G by sweeping in reverse topological order over the contracted graph (right side).

algorithm on each adjacent pair of blocks in the *quotient graph* where at least one of both is *active*. Initially all blocks are *active*. A block becomes *inactive* if none of its nodes move in a round. The algorithm terminates, if all blocks are *inactive*.

Fiduccia and Mattheyses [16] introduce a linear time two-way local search heuristic, called *FM* heuristic, which is fundamental for many graph partitioning algorithms. They define the gain $g(v)$ of a node $v \in V$ as the reduction of the cut metric when moving v from its current block to an other block. By maintaining the gains of the nodes in a special data structure, called *bucket queue*, they can find a maximum gain node in constant time. After moving a maximum gain node, they are also able to update the data structure in time equal to the number of adjacent nodes.

The local improvement algorithm (for *Active Block Scheduling*) can either be an *FM* local search or a flow-based approach or even a combination of both as proposed by Sanders and Schulz [42].

3.4. Hypergraph Partitioning

In this Section, we review how most hypergraph partitioners solve the *hypergraph partitioning problem*. The most successful approach is the *multilevel paradigm* [3, 5, 37] which we describe in Section 3.4.1. The results of this thesis is integrated into n -level hypergraph partitioner *KaHyPar*. Therefore, we give a brief overview of implementation details of this framework (see Section 3.4.2).

3.4.1. Multilevel Paradigm

The *multilevel paradigm* is a three phase algorithm to solve the *hypergraph partitioning problem* (see Figure 10). In the first stage, called *coarsening phase*, vertex matchings or clusterings are calculated which we contract. This process is repeated until a predefined number of hypernodes remains. The sequence of successively smaller hypergraphs is called *levels*. If the hypergraph H is small enough, we can use expensive algorithms to initially partition H into k blocks (*Initial Partitioning*). Afterwards, we *uncontract* each *level* in reverse order of *contraction* by projecting the partition to the next *level*. After *uncontraction* a *refinement* heuristic can be used to improve the quality of the current partition according to an objective function. The most commonly used *refinement* algorithm is the *FM* algorithm [16].

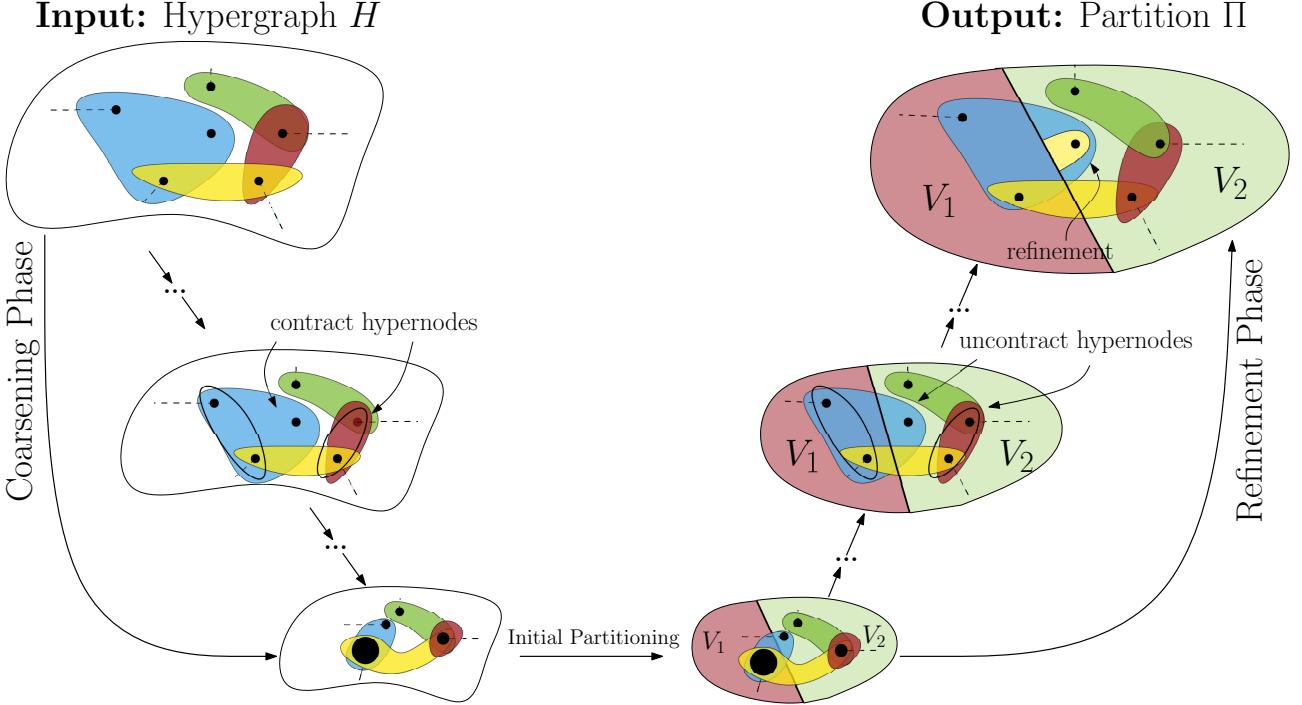


Figure 10: Multilevel Hypergraph Partitioning

3.4.2. n -Level Hypergraph Partitioning

KaHyPar is a multilevel hypergraph partitioner in its most extreme version, which removes only a single vertex in one *level* of the hierarchy. It seems to be the method of choice for optimizing cut- and the $(\lambda - 1)$ -metric unless speed is more important than quality [22]. The framework provides a *direct k-way* [1] and a *recursive bisection* mode, which recursively calculates bipartitions (with *multilevel paradigm*) until the hypergraph is divided into k blocks [43]. *KaHyPar* consists of four phases: *Preprocessing* and the three phases of the *multilevel paradigm*.

In the *preprocessing* step community structures of the hypergraph are detected. The hypergraph is transformed into a bipartite graph $G_*(H)$ (see Section 2.3) and a community detection algorithm is executed which optimizes *modularity* [19, 22]. During the *coarsening phase* contractions are restricted to vertices within the same community. The contraction partners are chosen according to the *heavy-edge* rating function $r(u, v) := \sum_{e \in I(u) \cap I(v)} \frac{\omega(e)}{|e|-1}$ [26]. The function prefers vertices which share a large number of heavy nets with small size. The contraction algorithm works in passes. At the beginning of each pass a random permutation of the vertices is generated and for each vertex u , we determine the contraction partner v according to the *heavy-edge* rating function [43]. A pass ends if each vertex is either considered as representative or contraction partner. The passes are repeated until only $t = 160k$ hypernodes remains. The *initial partitioning* uses the *recursive bisection* approach to calculate a k -way partition in combination with a portfolio of initial partitioning techniques [21]. In the *refinement phase*, a localized *FM* search is started [16], initialized with the current uncontracted vertices. The *local search* maintains k priority queues (PQ) for each block V_i exactly one [1]. A hypernode v contained in the i -th PQ with gain g means that moving vertex v to block V_i has gain g . After a move, the gains of all adjacent hypernodes are updated with a *delta-gain* update strategy [37]. The recalculation of all gain values at the beginning of a *FM* pass is one of the main bottlenecks of the algorithm [37]. Therefore, Schlag et al. [1, 43] introduce a *gain cache*, which prevents expensive recalculations of the corresponding gain function. The *gain cache* is maintained with *delta-gain* updates in the same way as the *PQs*. Further, the *local search* is stopped, when an

improvement during an *FM* pass becomes unlikely. This model is called *adaptive stopping rule* [1]. Sanders and Osipov [36] shows that it is unlikely that *local search* gives an improvement if $p > \frac{\sigma^2}{4\mu^2}$, where p is the number of moves in the current *FM* pass, μ is the average gain, and σ^2 the corresponding variance.

4. Hypergraph Flow Networks

In Section 3.2 we have shown how a hypergraph H can be transformed into a flow network $T_L(H)$ such that each minimum-weight (S, T) -cutset of H is a minimum-capacity (S, T) -cutset of $T_L(H)$ [30]. However, the resulting flow network has significantly more nodes and edges than the original hypergraph. The running time of a maximum (S, T) -flow algorithm depends heavily on the problem size. Therefore, different modeling approaches, which reduce the number of nodes and edges, can have a crucial impact on the running time of the flow algorithm.

We will present techniques to sparsify the flow network proposed by Lawler. First, we will show how *any subset $V' \subseteq V$ of hypernodes* could be removed from $T_L(H)$ (see Section 4.1). This approach minimizes the number of nodes, but in some cases, the number of edges can be significantly higher than in $T_L(H)$. The basic idea of this technique can still be applied to remove low degree hypernodes from the Lawler-Network *without* increasing the number of edges (see Section 4.2). Additionally, we show how every hyperedge e of size 2 can be removed by inserting an undirected flow edge between the corresponding nodes (see Section 4.3). Finally, we combine the two suggested approaches into a Hybrid-Network (see Section 4.4).

4.1. Removing Hypernodes via Clique-Expansion

In this Section, we show how all hypernodes of $T_L(H)$ can be removed such that a maximum (S, T) -flow on the new network induce a minimum-weight (S, T) -cutset on H . If a hypernode $v \in V$ occurs in an augmenting path P the previous node in the path must be a hyperedge, either e' or e'' . Further, for all $e \in I(v)$ the capacity $u_L(v, e')$ is ∞ . Therefore, if we push flow over a hypernode v , coming from a hyperedge, we can redirect the flow to any hyperedge node $e' \in I(v)$ during the whole maximum flow calculation because $u_L(v, e') = \infty$. The following lemma is central to our first sparsifying technique and is illustrated in Figure 11. Given a graph $G = (V, E)$ we define the two sets $in(u) := \{v \mid (v, u) \in E\}$ and $out(u) := \{v \mid (u, v) \in E\}$ with $u \in V$.

Lemma 4.1 (Shortcut Edges). *Let $G = (V, E, u)$ be a flow network and $u \in V$ a node where all incoming and outgoing edges have capacity equal to ∞ . Further, let $G(u) = (V \setminus \{u\}, E_u, u_u)$ be the flow network obtained by removing u and inserting a shortcut edge between each $v \in in(u)$ and $w \in out(u)$ with $u_u(v, w) = \infty$. If f is a maximum (S, T) -flow of G with $|f| < \infty$, then f is equal to a maximum (S, T) -flow f' of $G(u)$ with $u \notin S \cup T$.*

Proof. Let f be a maximum (S, T) -flow of G . We define a maximum (S, T) -flow f' of $G(u)$ as follows:

$$f'(v, w) = \begin{cases} \frac{f(v, u)f(u, w)}{\sum_{w \in out(u)} f(u, w)}, & \text{if } v \in in(u), w \in out(u) \\ f(v, w), & \text{otherwise} \end{cases} \quad (4.1)$$

f' is chosen in such a way that for all $v \in in(u) : \sum_{w \in out(u)} f'(v, w) = f(v, u)$ and for all $w \in out(u) : \sum_{v \in in(u)} f'(v, w) = f(u, w)$. Therefore, f' satisfies the flow conservation constraint and since all capacities are equal to ∞ , f' also satisfies the capacity constraint $\Rightarrow f'$ is a valid flow function. Further u is not contained in $S \cup T$ which implies that $|f| = |f'|$.

Let f' be a maximum (S, T) -flow of $G(u)$. We define a maximum (S, T) -flow f of G as follows:

$$\begin{aligned} f(u, w) &= \sum_{x \in in(u)} f'(x, w) \\ f(v, u) &= \sum_{x \in out(u)} f'(v, x) \\ f(x, y) &= f'(x, y) \end{aligned} \quad (4.2)$$

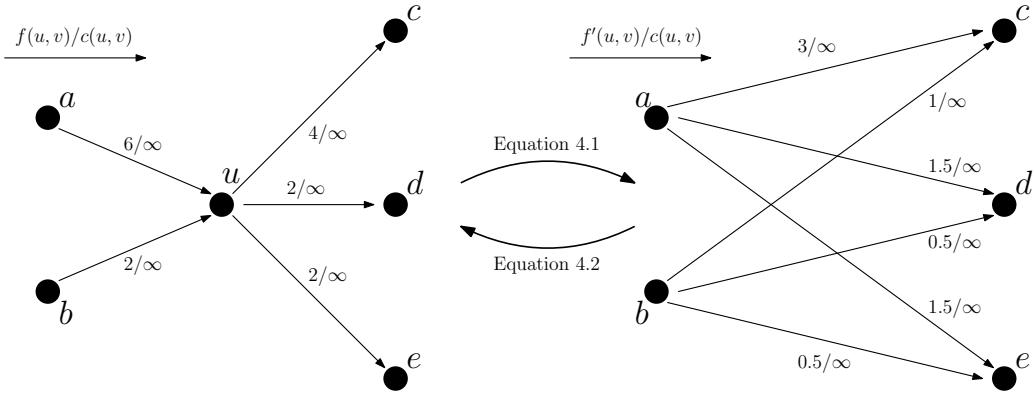


Figure 11: Illustration of Lemma 4.1 and Equation 4.1 and 4.2.

The amount of flow from each $v \in \text{in}(u)$ to each $w \in \text{out}(u)$ of flow function f' is redirected over u in f . Therefore, f is a valid flow function. Since $u \notin S \cup T$, it follows that $|f| = |f'|$. \square

In $T_L(H)$ all incoming and outgoing edges of a hypernode v have capacities equal to ∞ . For all $e \in I(v)$ there is an edge from v to e' and from e'' to v . Consequently, $\text{in}(v) = \bigcup_{e \in I(v)} e''$ and $\text{out}(v) = \bigcup_{e \in I(v)} e'$. Therefore, we can construct the following network with Lemma 4.1:

Definition 4.1. Let T_H be a transformation that converts a hypergraph $H = (V, E, c, \omega)$ into a flow network $T_H(H, V') = (V_H, E_H, u_H)$ with $V' \subseteq V$. $T_H(H, V')$ is defined as follows:

- (i) $V_H = V \setminus V' \bigcup_{e \in E} \{e', e''\}$
- (ii) $\forall v \in V'$ and $\forall e_1, e_2 \in I(v)$ with $e_1 \neq e_2$ we add a directed edge (e_1'', e_2') with capacity $u_H(e_1'', e_2') = \infty$ (Lemma 4.1).
- (iii) Let $H' = (V \setminus V', E', c, \omega)$ be the hypergraph with $E' = \{e \setminus V' \mid e \in E \wedge e \setminus V' \neq \emptyset\}$, then we add all edges of $T_L(H')$ to E_H with their corresponding capacities.

An example of the transformation is shown in Figure 12. We have to proof that a minimum-capacity (S, T) -cutset of $T_H(H, V')$ is equal with a minimum-weight (S, T) -cutset of H . However, we will use the following lemma in the correctness proof.

Lemma 4.2 (Source/Sink Node Removal). Let $G = (V, E, u)$ be a flow network and f a maximum (S, T) -flow of G with $|f| < \infty$. If $s \in S$ is a source node where all outgoing edges have infinite capacity and $t \in T$ is a sink node where all incoming edges have infinite capacity, then $|f|$ is equal with the amount of a maximum (S', T) -flow f_s of $G(s)$ and a maximum (S, T') -flow f_t of $G(t)$, where $S' = S \setminus \{s\} \cup \text{out}(s)$ and $T' = T \setminus \{t\} \cup \text{in}(t)$.

Proof. First we note, that the flow over an incoming edge of a source node $s \in S$ is zero. More formally, $\forall v \in \text{in}(s) : f(v, s) = 0$. Edmond and Karp [15] show that we can find a maximum (s, t) -flow if we augment in each step along a shortest path. Assume we find an augmenting path P which contains an edge (v, s) . We can obtain a shorter path if we split P after edge (v, s) and use the second part as augmenting path. Therefore, $f(v, s) = 0$. The same holds for all outgoing edges of a sink node. Consequently, we can remove all incoming resp. outgoing edges of a source resp. sink node.

In Section 2.2 we described how to solve a *multi-source multi sink* flow problem by adding a super source node a and super sink node b to the network and connect a with all sources $s' \in S$ and all sinks $t' \in T$ with b . $\forall s' \in S : u(a, s') = \infty$ and $\forall t' \in T : u(t', b) = \infty$. With Lemma 4.1 follows, that we can remove s from G and insert a directed edge from a to each $v \in \text{out}(s)$ (equal to $G(s)$) and $|f| = |f_s|$. The new flow problem corresponds to the *multi-source multi-sink* problem with S' and T as source and sink set. The proof for $G(t)$ is equivalent. \square

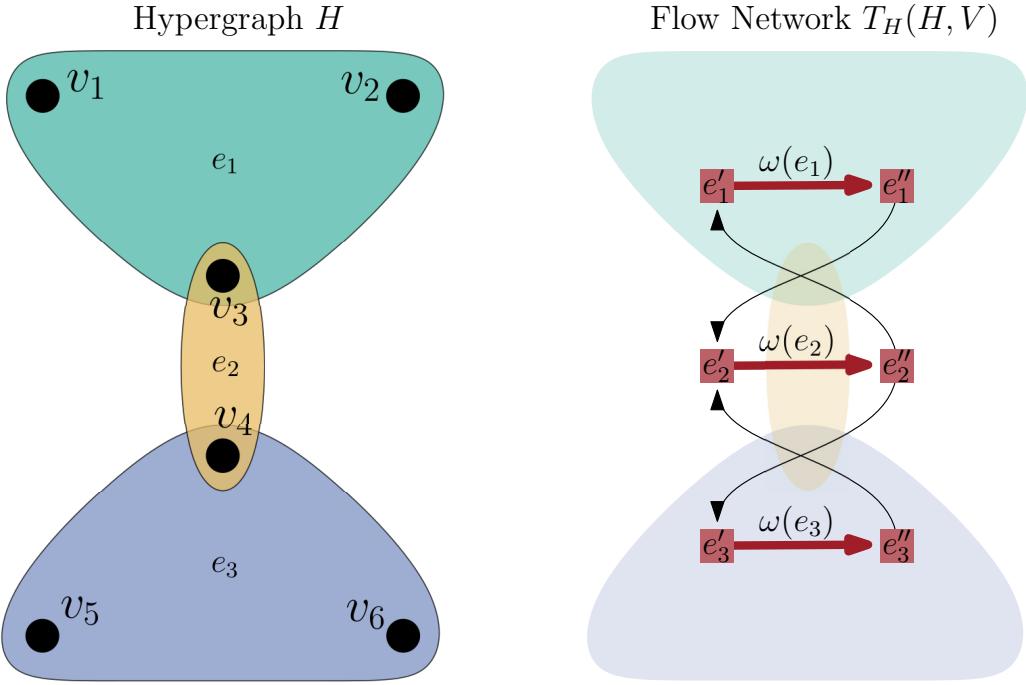


Figure 12: Transformation of a hypergraph H into an equivalent flow network $T_H(H, V)$ by removing all hypernodes of $T_L(H)$. Note, capacity of the black edges in the flow network is ∞ .

As a consequence of this lemma, we can remove a source hypernode $v \in S$ of $T_L(H)$ and instead add all incoming hyperedge nodes $e' \in I(v)$ as sources to the flow problem. Because for all incoming resp. outgoing edges of vertices v of $T_L(H)$ the capacity is ∞ .

Theorem 4.1. *A minimum-weight (S, T) -cutset of a hypergraph $H = (V, E, c, \omega)$ (with $S, T \subseteq V, S \cap T = \emptyset$) is equivalent with a minimum-capacity (S', T') -cutset of the flow network $T_H(H, V')$ ($V' \subseteq V$) with $S' = S \setminus V' \cup \bigcup_{e \in N(V' \cap S)} \{e'\}$ and $T' = T \setminus V' \cup \bigcup_{e \in N(V' \cap T)} \{e''\}$.*

Proof. Applying Lemma 4.1 and 4.2 on all nodes $v \in V'$ of flow network $T_L(H)$ yields network $T_H(H, V')$ with S' and T' as source and sink sets. A maximum (S, T) -flow f_L of $T_L(H)$ is then equal with a maximum (S', T') -flow f_H of $T_H(H, V')$. Since $|f_L| < \infty$, only edges between hyperedge nodes are contained in a minimum-capacity (S, T) -cutset of $T_L(H)$. Since $|f_L| = |f_H|$, the same holds for a minimum-capacity (S', T') -cutset of $T_H(H, V')$, which is equal with a minimum-weight (S, T) -cutset of H .

□

Consequently, we can find a minimum-weight (S, T) -cutset of H by calculating a minimum-capacity (S', T') -cutset of $T_H(H, V')$. Finally, we have to find the corresponding minimum-weight (S, T) -bipartition. In $T_L(H)$ all hypernodes reachable from source nodes in the residual graph are part of the first and all not reachable are part of the second block of the bipartition. Since we removed all hypernodes $v \in V'$ in our new network, we have to reconstruct the bipartition using the following lemma.

Lemma 4.3 (Reachability of Hypernodes). *Let f be a maximum (S, T) -flow of $T_L(H)$. If a hypernode $v \notin S$ is reachable from a node $s \in S$ in the residual graph of $T_L(H)$, then there must exist at least one net $e \in I(v)$ where e'' is reachable from s in the residual graph of $T_L(H)$.*

Proof. Let A be the set of all nodes reachable from the source nodes S in the residual graph of $T_L(H)$. Assume, if $v \in A$, then $\forall e \in I(v)$ the *outgoing hyperedge node* e'' is not contained in A which implies that all edges (v, e'') are not contained in the residual graph of $T_L(H)$. More formally, $\forall e \in I(v) : r_f(v, e'') = 0$. Otherwise, e'' would be in A because $v \in A$. Since $r_f(v, e'') = f(e'', v) = 0$, there is no flow entering node v and due to the conservation of flow constraint there cannot be any flow leaving node v . Therefore, there is no path from any $s \in S$ to v over a node e' , because $\forall e \in I(v) : r_f(e', v) = f(v, e') = 0$ and no path over e'' because $\forall e \in I(v) : e'' \notin A$. Therefore, v is not reachable from any $s \in S$ which is a contradiction to the assumption that $v \in A$. \square

Lemma 4.3 gives us an alternative construction for the minimum-weight (S, T) -bipartition of H for both networks $T_L(H)$ and $T_H(H, V')$. Regardless of the flow network, we can calculate a maximum flow on it and define the set E'' , which contains all *outgoing hyperedge nodes* e'' *reachable* from a source node $s \in S$ in the *residual graph* of the flow network. Further, $(A := \bigcup_{e \in E''} e, V \setminus A)$ is a minimum-weight (S, T) -bipartition of H .

4.2. Low-Degree Hypernodes

The resulting flow network $T_H(H, V)$ proposed in Section 4.1 has significantly fewer nodes than the network $T_L(H)$ proposed by Lawler. On the other hand, the number of edges could be much larger.

Consider a hypernode $v \in V$. We replace v in $T_L(H)$ with a biclique between all e'' and e' which are incident to v . The number of edges added to $T_H(H, V)$ depends on the degree of v . Each vertex $v \in V$ induces $d(v)(d(v) - 1)$ edges in $T_H(H, V)$. In $T_L(H)$, a hypernode adds $2d(v)$ edges to the network and add one additional node. A simple observation is that for all hypernodes with $d(v) \leq 3$ the inequality $d(v)(d(v) - 1) \leq 2d(v)$ holds. Removing such low degree hypernodes not only reduces the number of nodes, but also the number of edges.

Let $V_d(n) = \{v \in V \mid d(v) \leq n\}$ be the set of all hypernodes with degree smaller or equal n . Then our suggested flow network is $T_H(H, V_d(3))$.

4.3. Removing Graph Hyperedges

If we want to find a minimum-weight (S, T) -cutset of a graph $G = (V, E, \omega)$, we do not have to transform G into an equivalent flow network. We can directly operate on the graph with capacities $u(e) = \omega(e)$ for all $e \in E$ [17]. Hypergraphs are a generalization of graphs, where an edge can consist of more than two nodes. However, a hyperedge e of size 2 can still be interpreted as a graph edge. Instead of modeling those edges as described by Lawler [30] (see hyperedge e_2 in Figure 6), we can add an undirected flow edge between $v_1, v_2 \in e$ (with $v_1 \neq v_2$) with capacity $u(\{v_1, v_2\}) = \omega(e)$. In the following, we will proof the opposite. We will show that each undirected graph can be modeled as a directed graph with the same *min-cut* properties. The transformation used in the proof of an undirected to an directed edge will have the same structure as a hyperedge of size two in the Lawler-Network. As a consequence, if we define the network where each hyperedge of size two is modeled with an undirected flow edge, we can use the following lemma to show that both networks have the same value of a maximum (S, T) -flow.

Lemma 4.4 (Transformation of Undirected to Directed Networks). *Let $G = (V, E, u)$ be an undirected flow network with capacity function $u : E \rightarrow \mathbb{N}_+$. G can be transformed into a directed graph G' such that the value of a maximum (s, t) -flow f of G is equal with the value of a maximum (s, t) -flow f' of G' . More formally, $|f| = |f'|$.*

Proof. Assume $\forall e \in E : u(e) = 1$. According to Menger's Theorem [34], a maximum (s, t) -flow is then equal with the maximum number of edge-disjoint paths between s and t in a directed graph. This theorem can also be proven for undirected graphs if we replace each undirected edge $e = \{u, v\}$ by five directed edges $(v, x'), (w, x'), (x', x''), (x'', v), (x'', w)$ (see Figure 13) [34]. Obviously, we can map each set of edge-disjoint paths from s to t from G' to G and vice versa. Therefore, the maximum number of edge-disjoint paths from s to t in G' is then the same as G and therefore, $|f| = |f'|$.

Consider the general case where $\forall e \in E : u(e) \in \mathbb{N}_+$. We can transform the weighted undirected graph G into an unweighted directed multigraph by replacing each undirected edge $e = \{u, v\}$ with $u(e)$ undirected edges of weight 1 (see Figure 13). Afterwards, we can use the transformation to an unweighted directed multigraph the same way as before. Again, we can apply Menger's Theorem to show that $|f| = |f'|$. Newman [35] showed that there is an one-to-one correspondence between a maximum (s, t) -flow of an unweighted multigraph and its corresponding weighted graph where the weight of each edge (u, v) is the number of parallel edges between u and v of the multigraph. \square

As a consequence of the construction of the proof of Lemma 4.4 the weighted directed graph illustrated on the right side of Figure 13 can be transformed into a single undirected edge with weight $u(\{u, v\}) = u(x', x'')$. Each hyperedge e with $|e| = 2$ has exactly this structure in $T_L(H)$. Therefore, we can construct the following network:

Definition 4.2. Let T_G be a transformation that converts a hypergraph $H = (V, E, c, \omega)$ into a flow network $T_G(H) = (V_G, E_G, u_G)$. $T_G(H)$ is defined as follows:

- (i) $V_G = V \cup \bigcup_{e \in E : |e|=2} \{e', e''\}$
- (ii) $\forall e \in E$ with $|e| = 2$ and $v_1, v_2 \in e$ ($v_1 \neq v_2$) we add two directed edges (v_1, v_2) and (v_2, v_1) to E_G with capacity $u_G(v_1, v_2) = \omega(e)$ and $u_G(v_2, v_1) = \omega(e)$
- (iii) Let $H' = (V, E', c, \omega)$ be the hypergraph with $E' = \{e \mid e \in E \wedge |e| \neq 2\}$, then we add all edges of $T_L(H')$ to E_G with their corresponding capacities.

An example of transformation $T_G(H)$ is shown in Figure 14. A hyperedge e of size 2 consists exactly of 4 nodes and 5 edges in $T_L(H)$ (see Figure 6). The same hyperedge induces 2 nodes and 2 edges in $T_G(H)$ (see Figure 15).

Theorem 4.2. A minimum-weight (S, T) -cutset of a hypergraph $H = (V, E, c, \omega)$ (with $S, T \subseteq V, S \cap T = \emptyset$) is equal with a minimum-capacity (S, T) -cutset of the flow network $T_G(H) = (V_G, E_G, u_G)$.

Proof. Consider a hyperedge e with $|e| = 2$ and $u, v \in e$ ($u \neq v$). The capacity of $(u, e'), (v, e')$, (e'', u) and (e'', v) is infinity in flow network $T_L(H)$. Before we can apply Lemma 4.4 on all hyperedges e with $|e| = 2$, we have to show how to handle the infinite capacity edges. The flow leaving e' is restricted by $u(e', e'') = \omega(e)$. Therefore, the flow entering e' is restricted by $f(u, e') + f(v, e') \leq u(e', e'') = \omega(e)$. Consequently, $f(u, e') \leq \omega(e)$ and $f(v, e') \leq \omega(e)$. The same holds for $f(e'', u)$ and $f(e'', v)$. Therefore, we can replace each infinite capacity of an edge entering e' or leaving e'' with $\omega(e)$ without changing the value of a maximum (S, T) -flow. We call the capacity adapted network $T_{L'}(H)$.

Applying the transformation of Lemma 4.4 on each undirected edge of $T_G(H)$ results in flow network $T_{L'}(H)$. It follows, that a maximum (S, T) -flow of $T_G(H)$ is equal with a maximum (S, T) -flow of $T_{L'}(H)$ and $T_L(H)$. Consequently, a minimum-capacity (S, T) -cutset of $T_G(H)$ is equal with a minimum-weight (S, T) -cutset of H . \square

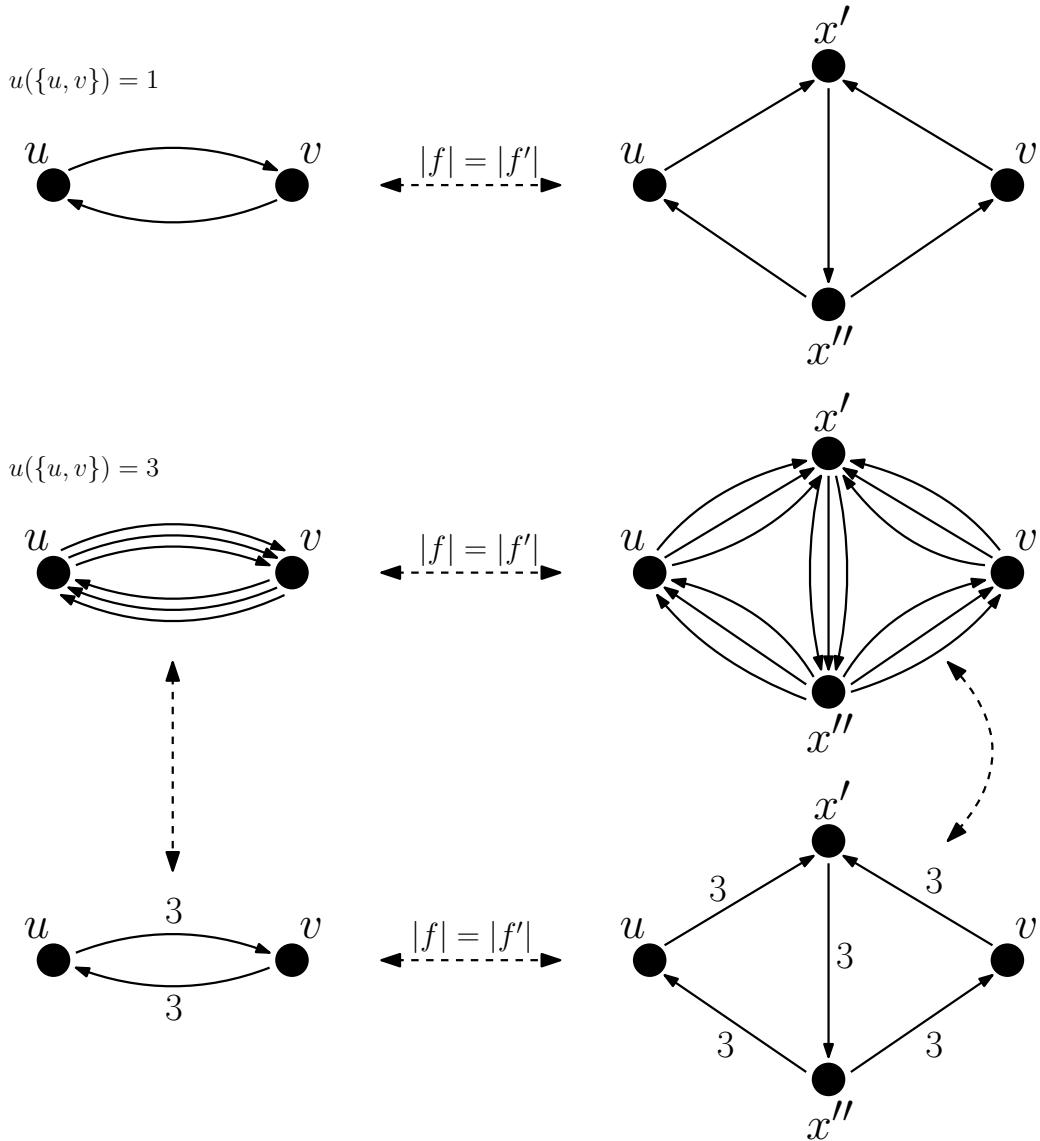


Figure 13: Illustration of the transformation of an unweighted or weighted undirected graph into an unweighted or weighted directed graph. The equivalence of a maximum (s, t) -flow of a unweighted multigraph and their corresponding weighted graph is a result of a work by Newman [35].

A minimum-weight (S, T) -cutset of H can also be calculated with $T_G(H)$. Each edge (v_1, v_2) with $v_1, v_2 \in V$ of the minimum-capacity (S, T) -cutset of $T_G(H)$ can be mapped to their corresponding hyperedge. Since there exists a one-one correspondence between the hypernodes of $T_L(H)$ and $T_G(H)$ the corresponding bipartition are all hypernodes *reachable* from all nodes in S and all not *reachable* from S in the *residual graph* of $T_G(H)$.

4.4. Combining Techniques

The density of a hypergraph $H = (V, E)$ is defined as follows:

$$d := \frac{\overline{d(v)}}{|e|} = \frac{|P|/|V|}{|P|/|E|} = \frac{|E|}{|V|}$$

where $\overline{d(v)}$ is the average hypernode degree, $\overline{|e|}$ is the average hyperedge size and $|P|$ are the number of pins. Many real world benchmark instances have either a low or high density. For

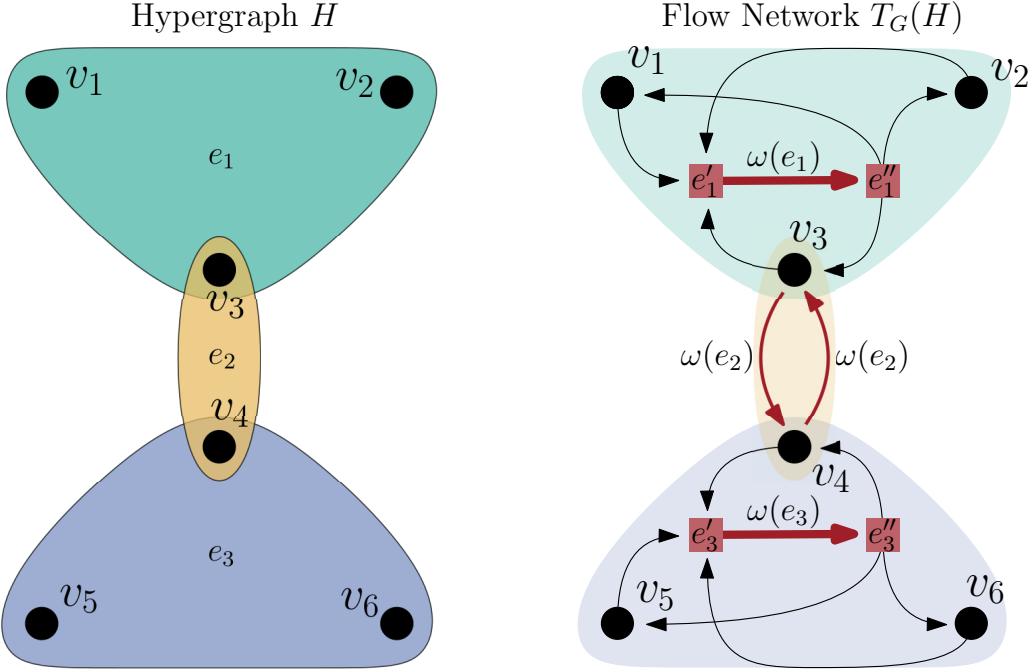


Figure 14: Transformation of a hypergraph into an equivalent flow network by inserting an undirected edge with capacity $\omega(e)$ for each hyperedge of size 2. Note, capacity of the black edges in the flow network is ∞ .

an example, consider the statistical summary of our benchmark set in Table 7. Hypergraphs with a large density have usually a average hypernode degree significantly greater than the average hyperedge size. Whereas, the opposite behavior can be observed on instances with a low density. High density hypergraphs often have a structure similiar to a graph and low density hypergraphs often have large hyperedges with low degree hypernodes. Of course, there exists instances which are counterexamples to this observation, but for a large majority of real world benchmarks we often find the described behavior.

Currently, we have two different modeling approaches which either perform better on hypergraphs with many low degree hypernodes or small hyperedges. Taking our observation from real-world instances into account means that $T_G(H)$ performs significantly better on high density hypergraphs and $T_H(H, V_d(3))$ on low density hypergraphs. It would be preferable to combine the two approaches into one network which performs best on most instances.

Definition 4.3. Let T_{Hybrid} be a transformation that converts a hypergraph $H = (V, E, c, \omega)$ into a flow network $T_{\text{Hybrid}}(H, V') = (V_{\text{Hybrid}}, E_{\text{Hybrid}}, u_{\text{Hybrid}})$, where $V' = \{v \in V_d(3) \mid \forall e \in I(v) : |e| \neq 2\}$. $T_{\text{Hybrid}}(H, V')$ is defined as follows:

- (i) $V_{\text{Hybrid}} = V \setminus V' \cup \bigcup_{\substack{e \in E \\ |e| \neq 2}} \{e', e''\}$
- (ii) $\forall v \in V'$ we add a directed edge (e''_1, e'_2) , $\forall e_1, e_2 \in I(v)$ ($e_1 \neq e_2$) with capacity $u_{\text{Hybrid}}(e''_1, e'_2) = \infty$ (Lemma 4.1).
- (iii) $\forall e \in E$ with $|e| = 2$ and $v_1, v_2 \in e$ we add two directed edges (v_1, v_2) and (v_2, v_1) with capacity $u_{\text{Hybrid}}(v_1, v_2) = \omega(e)$ and $u_{\text{Hybrid}}(v_2, v_1) = \omega(e)$ (Lemma 4.4)
- (iv) $\forall e \in E$ with $|e| \neq 2$ we add a directed edge (e', e'') with capacity $u_{\text{Hybrid}}(e', e'') = \omega(e)$ (same as in $T_L(H)$).
- (v) $\forall v \in V \setminus V'$ we add for each incident hyperedge $e \in I(v)$ with $|e| \neq 2$ two directed edges (v, e') and (e'', v) with capacity $u_{\text{Hybrid}}(v, e') = u_{\text{Hybrid}}(e'', v) := \infty$ (same as in $T_L(H)$).

Figure 15 summarizes all explained transformations of this section. We can prove the correctness of $T_{\text{Hybrid}}(H, V')$ with Lemma 4.1, 4.2 and 4.4 as used in the proof of Theorem 4.1 and 4.2. A minimum-weight (S, T) -cutset of H is equal with a minimum-capacity (S', T') -cutset of $T_{\text{Hybrid}}(H, V')$.

Per definition of $T_{\text{Hybrid}}(H, V')$ we prefer hyperedge removal over hypernode removal. If a hypernode has a degree smaller than or equal to 3, we only remove it, if there is no hyperedge $e \in I(v)$ with $|e| = 2$. The reason for this is that hyperedge removal always removes more nodes and edges than hypernode removal.

The minimum-weight (S, T) -cutset of H can be calculated using the technique described in Section 4.3. Let $(A, V \setminus A)$ be the corresponding bipartition. A is the union of all reachable hypernodes from S' and the union of all reachable *outgoing hyperedge nodes* e'' from S' (see Section 4.1 and Lemma 4.3).

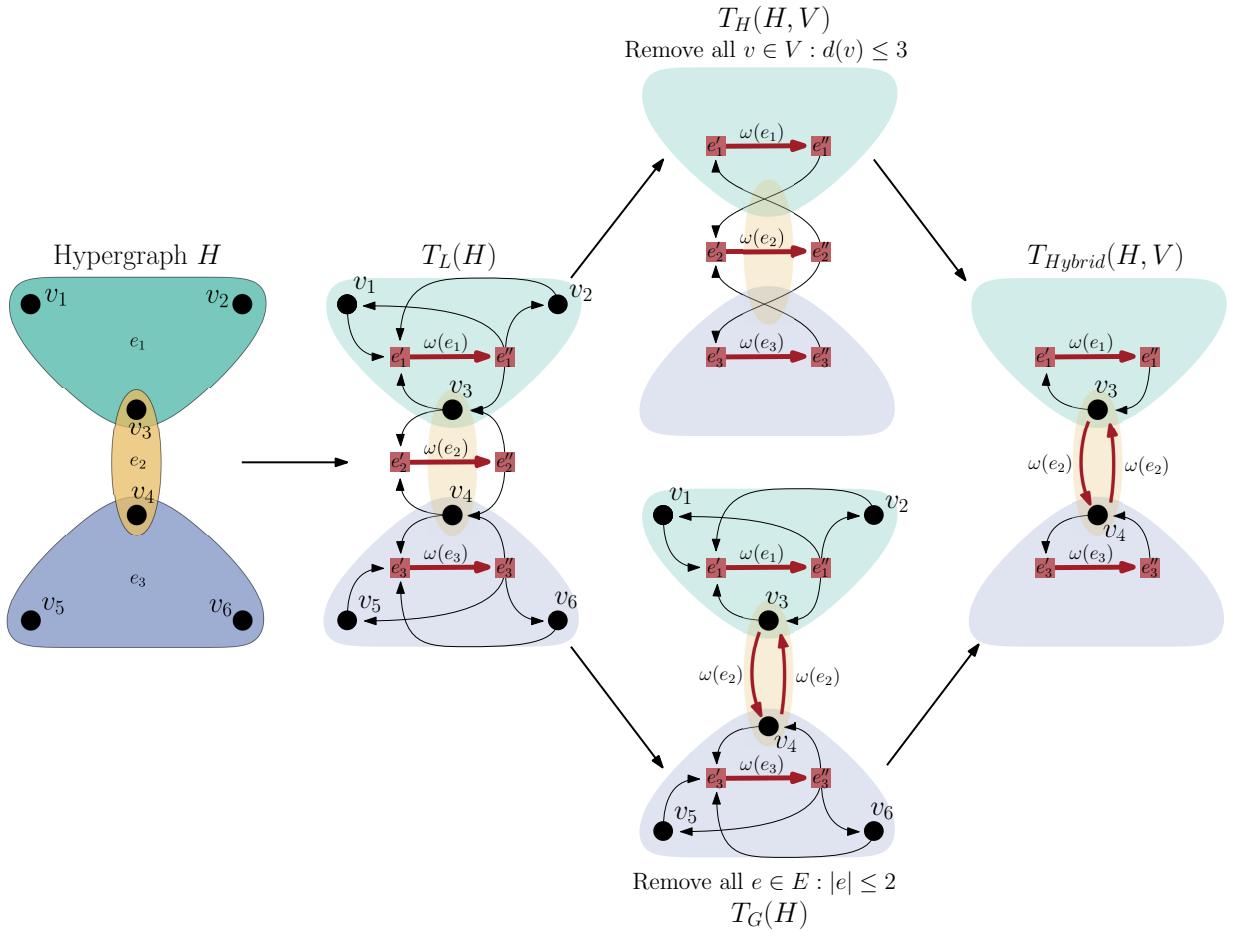


Figure 15: Illustration of all presented techniques to sparsify the flow network of a hypergraph. Transformation from $T_L(H)$ to $T_H(H, V)$ follows with Lemma 4.1. Transformation from $T_L(H)$ to $T_G(H)$ follows with Lemma 4.4.

5. Max-Flow-Min-Cut Refinement Framework

We will now give a detailed overview of our flow-based refinement framework. The main idea is to extract a subhypergraph $H_{V'}$ out of a hypergraph H , which is already partitioned into k blocks. V' is chosen in such way that it is a subset of two adjacent blocks V_i and V_j . We will show how to configure the sources S and sinks T of the corresponding flow network such that a minimum (S, T) -bipartition of $H_{V'}$ improves the connectivity metric of H (see Section 5.1). Further, we describe how the ideas of Sanders and Schulz [42] can be adapted to work in an n -level hypergraph partitioner, called *KaHyPar*.

5.1. Source and Sink Configuration

Let $H = (V, E, c, \omega)$ be a hypergraph and B_1 be a bipartition of H . In the following, we show how to configure the source set S and sink set T of the flow network $T_L(H_{V'})$ of a subhypergraph $H_{V'}$ induced by $V' \subseteq V$. The goal is to improve bipartition B_1 of H with a maximum (S, T) -flow calculation on $T_L(H_{V'})$ (with f as maximum flow) such that after inserting the minimum (S, T) -bipartition of $H_{V'}$ on H the resulting bipartition B_2 of H has a cut less than or equal to the cut of B_1 . Let $E_{\text{cut}}(V', B) := \{e \in E \mid \lambda(e, B) > 1 \wedge e \cap V' \neq \emptyset\}$ be the set of all cut hyperedges of bipartition B of H which are partially or fully contained in $H_{V'}$. We define the cut of subhypergraph $H_{V'}$ related to a bipartition B of H as follows:

$$\omega_{H_{V'}}(B) := \sum_{e \in E_{\text{cut}}(V', B)} \omega(e)$$

Note, the cut $\omega_{H_{V'}}$ is defined over the cut nets of H . A cut hyperedge e of H is not necessarily a cut hyperedge of $H_{V'}$. If $e = \{v_1, v_2\}$ is a hyperedge with $v_1 \in V_1$ and $v_2 \in V_2$ and $v_1 \in V'$ and $v_2 \notin V'$, then e is cut in H , but not in $H_{V'}$, because v_2 is removed from e of $H_{V'}$ per definition. However, the reason that we still define e as cut hyperedge of $H_{V'}$ has to do with our problem statement, which we will define as follows:

Problem 5.1. *How do we have to define the source set S and sink set T for a subhypergraph $H_{V'}$ (with $V' \subseteq V$) and a bipartition B_1 such that after a maximum (S, T) -flow calculation (with f as maximum flow) the resulting minimum (S, T) -bipartition B_2 of H satisfy the following conditions:*

- (i) $\omega_H(B_2) \leq \omega_H(B_1)$
- (ii) $\Delta_H := \omega_H(B_1) - \omega_H(B_2) = \omega_{H_{V'}}(B_1) - |f| =: \Delta_{H_{V'}}$

The first condition ensures that a maximum (S, T) -flow calculation on $T_L(H_{V'})$ never decrease the cut of H . The existence of the second condition has practical reasons. First, we can simply update the cut metric via $\omega_H(B_2) = \omega_H(B_1) - \Delta_{H_{V'}}$, instead of summing up the weight of all cut hyperedges. Since we have to setup the subhypergraph $H_{V'}$ before each maximum flow computation, we can implicitly calculate $\omega_{H_{V'}}(B_1)$. Therefore, the cut metric can be updated after a *Max-Flow-Min-Cut* computation in constant time instead of $\mathcal{O}(|E|)$. On the other hand, we can assert the correctness of our maximum flow algorithm. If $\Delta_H \neq \Delta_{H_{V'}}$, then with high probability our flow algorithm is incorrect. The reason why we define $\omega_{H_{V'}}(B)$ over the cut hyperedges of H is that the equality

$$\Delta_H := \omega_H(B_1) - \omega_H(B_2) = \omega_{H_{V'}}(B_1) - \omega_{H_{V'}}(B_2)$$

holds only if we use the adapted defintion. Further, if we can show that $|f| = \omega_{H_{V'}}(B_2)$, we simultaneously show that our source and sink set modeling approach satisfies condition (ii)

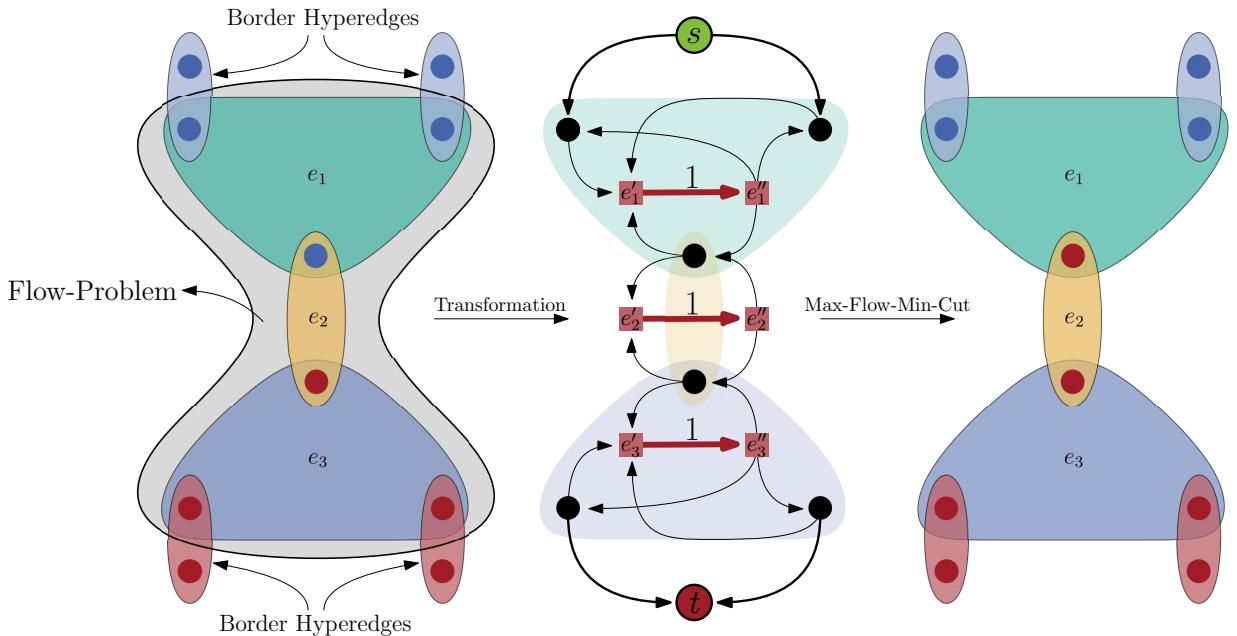


Figure 16: Non-cut *border hyperedges* of H and $H_{V'}$ induce source and sink hypernodes in the flow problem.

$$\Delta_H = \Delta_{H_{V'}}.$$

We will now present a solution for our problem statement. First, we show how S and T can be chosen to satisfy condition (i). Afterwards, we extend S and T with additional nodes to fulfil condition (ii).

Let $V' \subseteq V$ and $\delta B = \{e \in E \mid \exists u, v \in e : u \in V' \wedge v \notin V'\}$ be the set of all *border hyperedges* (see Figure 8). Further, we divide δB into two disjoint subsets:

- (i) Non-Cut hyperedges $e \in \delta B$ of H : $\delta B_1 = \{e \in \delta B \mid e \subseteq V_1 \vee e \subseteq V_2\}$
- (ii) Cut hyperedges $e \in \delta B$ of H : $\delta B_2 = \delta B \setminus \delta B_1$

For a bipartition (V_1, V_2) of H , we say $v \in V_1$ is a source node of the flow network $T_L(H_{V'})$, if there exists a hyperedge $e \in \delta B_1$ containing v . More formal:

$$S_1 = \{s \in V' \cap V_1 \mid \exists e \in \delta B_1 : s \in e\} \quad (5.1)$$

$$T_1 = \{t \in V' \cap V_2 \mid \exists e \in \delta B_1 : t \in e\} \quad (5.2)$$

An example of a *Max-Flow-Min-Cut* computation of $H_{V'}$ with S and T as source and sink set is illustrated in Figure 16.

Lemma 5.1. *Let B_1 be a bipartition of H and $T_L(H_{V'})$ the flow network of subhypergraph $H_{V'}$ with S and T as defined in Equation 5.1 and 5.2 (with $V' \subseteq V$). If B_2 is a bipartition obtained by a maximum (S, T) -flow computation on $T_L(H_{V'})$, then the inequality $\omega_H(B_2) \leq \omega_H(B_1)$ holds.*

Proof. A maximum (S, T) -flow computation on $T_L(H_{V'})$ yields a minimum (S, T) -cutset on $H_{V'}$ [17]. Thus, for all hyperedges $e \notin \delta B$ (fully contained in $H_{V'}$) which are cut in B_2 , the sum of their weight must be less or equal than the sum of all cut hyperedges $e \notin \delta B$ of bipartition B_1 . We have to show that a non-cut hyperedge $e \in \delta B_1$ of $B_1 = (V_1, V_2)$ cannot become a cut hyperedge of $B_2 = (V'_1, V'_2)$. Let $e \in \delta B_1$ be such a hyperedge. e must be either a subset of V_1 or V_2 , otherwise e is a cut hyperedge. Let $e \subseteq V_1$, then $e \cap V' \subseteq S$ (see Equation 5.1). Defining a node $s \in S$ as source node means that it cannot change its block after a *Max-Flow-Min-Cut*

computation. Therefore, $e \subseteq V_1$ and $e \subseteq V'_1 \Rightarrow e$ is a non-cut hyperedge of B_2 . The proof for $e \subseteq V_2$ is equivalent $\Rightarrow \omega_H(B_2) \leq \omega_H(B_1)$. \square

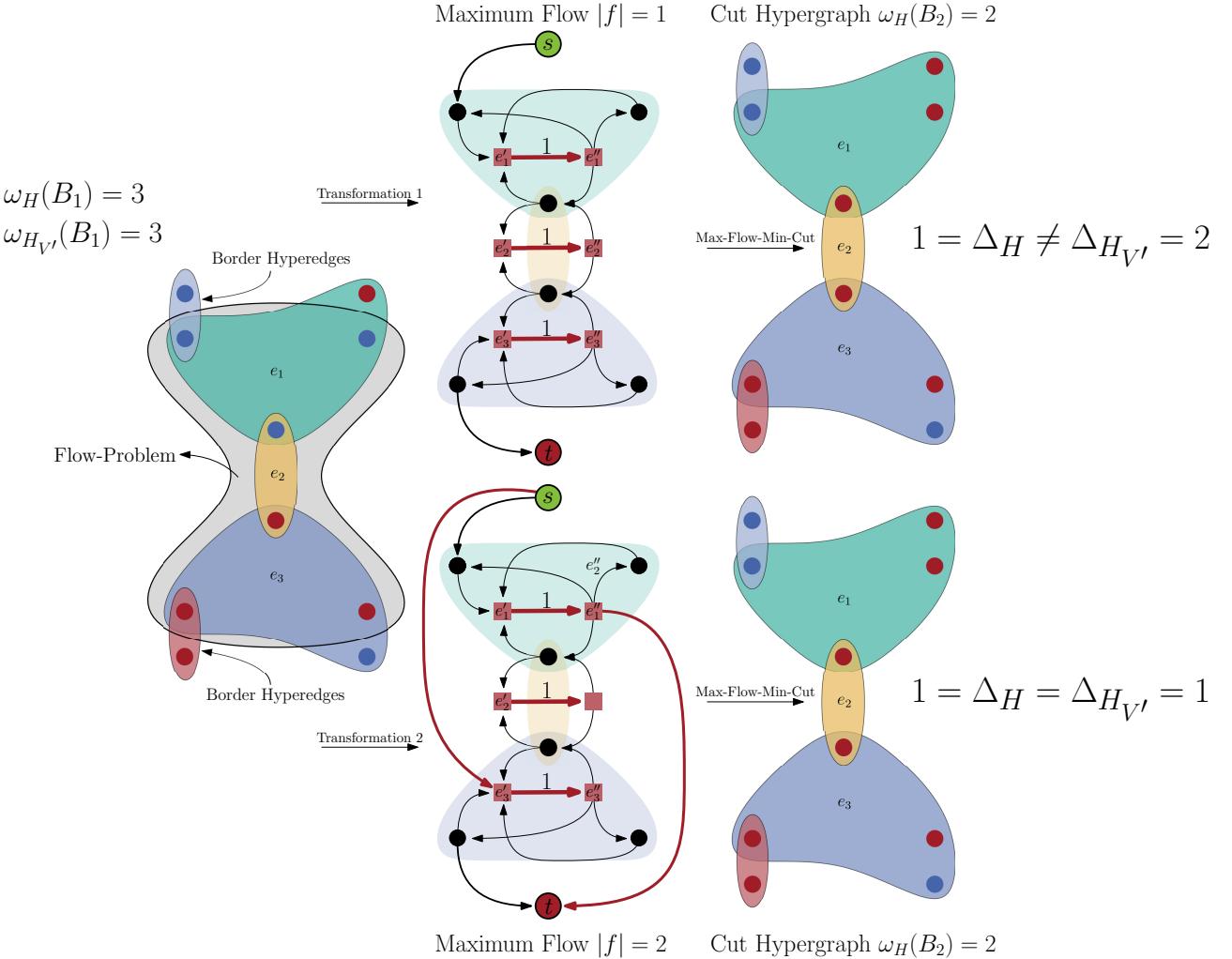


Figure 17: In this example e_1 and e_3 are cut hyperedges of the hypergraph, but non-cut nets of subhypergraph $H_{V'}$. Modeling the *outgoing* resp. *incoming* hyperedge node of e_1 resp. e_2 as sink resp. source ensures that $\Delta_H = \Delta_{H_{V'}}$.

In the next step, we will show how S and T can be extended to satisfy condition (ii) of Problem 5.1. Currently, $|f| \leq \omega_{H_{V'}}(B_2)$ (without a prove). Obviously, some nodes are missing in S and T . Consider Figure 17 to understand which nodes are missing. Transformation 1 illustrates our current modeling approach defined in Equation 5.1 and 5.2. The maximum flow on this network is $|f| = 1$, but the resulting minimum (S, T) -bipartition B_2 induce a cut of $\omega_H(B_2) = 2$. This implies that $\Delta_H = 3 - 2 \neq 3 - 1 = \Delta_{H_{V'}}$. The hyperedges e_1 and e_3 are cut nets of H , but non-cut hyperedges of $H_{V'}$. Therefore, B_1 induce a cut of 1 on $H_{V'}$ if we define the cut $\omega_{H_{V'}}(B_2)$ over the cut hyperedges of $H_{V'}$ instead of the cut hyperedges of H . In our example, we can remove e_2 from cut, but e_1 becomes a cut hyperedge of $H_{V'}$. Therefore, the value of the cut of $H_{V'}$ does not change, but the cut of H does. e_1 is already a cut hyperedge of H and B_2 removes e_2 from the cut of H . Therefore, $\Delta_H = 1$. However, we defined $\omega_{H_{V'}}(B_2)$ over the cut hyperedges of H and currently, we have $|f| = 1 \neq 2 = \omega_{H_{V'}}(B_2)$.

Transformation 2 illustrates the adapted modeling approach for cut hyperedges of H . For each hyperedge $e \in \delta B_2$ with $e \setminus V' \cap V_1 \neq \emptyset$, we add the *incoming* hyperedge node e' to S . More

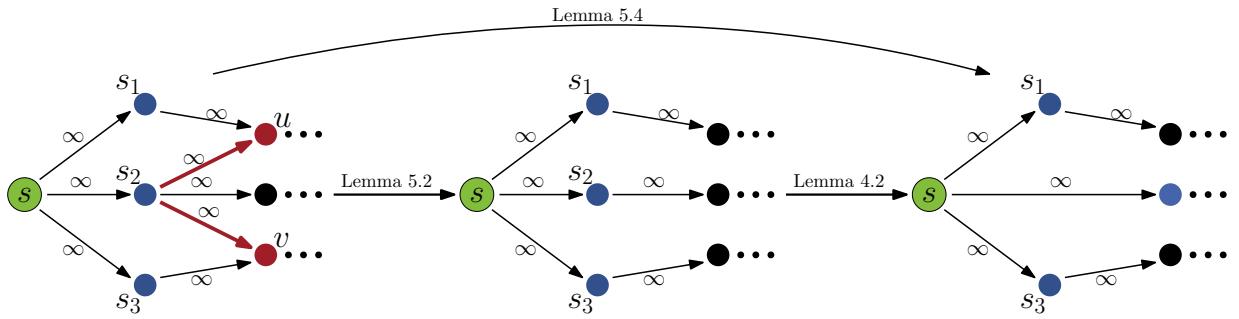


Figure 18: Illustration of the proof technique used in Lemma 5.4. The green node s is the super source of the flow problem. The blue nodes are source nodes of the corresponding *multi-source multi-sink* flow problem. The red nodes u and v are contained in $\mathcal{R}(s_2)$. Therefore, edges (s_2, u) and (s_2, v) are removable.

formal:

$$S = S_1 \cup \{e' \in \delta B_2 \mid e \setminus V' \cap V_1 \neq \emptyset\} \quad (5.3)$$

$$T = T_1 \cup \{e'' \in \delta B_2 \mid e \setminus V' \cap V_2 \neq \emptyset\} \quad (5.4)$$

Definition 5.1 (Extension of a Subhypergraph). *We define the extension $Ex(H_{V'})$ of a subhypergraph $H_{V'}$ such that each hyperedge of $H = (V, E, c, \omega)$ which is partially contained in $H_{V'}$ is fully contained in $Ex(H_{V'})$. More formally, $Ex(H_{V'}) = (V' \cup V'', E', c, \omega)$ with $V'' = \bigcup_{e \in \delta B} e \setminus V'$ and $E' = \{e \in E \mid e \subseteq (V' \cup V'')\}$.*

We have to show that for a maximum (S, T) -flow f of $T_L(H_{V'})$ holds $|f| = \omega_{H_{V'}}(B_2)$. The idea of the proof is to use the extension $Ex(H_{V'})$ of $H_{V'}$ and add all $V'' \cap V_1$ to S_1 (see Equation 5.1) and all $V'' \cap V_2$ to T_1 (see Equation 5.2). We can show that for a maximum (S_1, T_1) -flow f' of $T_L(Ex(H_{V'}))$ holds that $|f'| = \omega_{H_{V'}}(B_2)$.

Afterwards, we use a technique to remove all $v \in V''$ of $T_L(Ex(H_{V'}))$ and show that the resulting flow network is $T_L(H_{V'})$ with S and T as source and sink set as defined in Equation 5.3 and 5.4. Moreover, for a maximum (S, T) -flow f then holds that $|f| = |f'| = \omega_{H_{V'}}(B_2)$.

Because of the complexity of the proof, we will introduce lemmas in the following which will simplify the proof of the main theorem. Consider Figure 18 if you need an illustration for the following lemmas.

Lemma 5.2 (Source Edge Removal). *Let f be a maximum (S, T) -flow of $G = (V, E, u)$. If there exists two edges (s_1, v) and (s_2, v) with infinite capacity ($s_1, s_2 \in S$) we can either remove (s_1, v) or (s_2, v) from G without changing the amount of a maximum (S, T) -flow.*

Proof. Let $P = (s_1, v, \dots)$ be an augmenting path of G . Replacing s_1 in P with s_2 yields an augmenting path P' of same length. The operation is valid because $u(s_1, v) = u(s_2, v) = \infty$. If we execute Edmond and Karp's maximum flow algorithm we can map each augmenting path P to P' and ensure that for a maximum (S, T) -flow f follows that $f(s_1, v) = 0$. Consequently, there exists maximum (S, T) -flows where either $f(s_1, v) = 0$ or $f(s_2, v) = 0$. Therefore, we can remove either (s_1, v) or (s_2, v) without changing the amount of a maximum (S, T) -flow. \square

Lemma 5.3 (Sink Edge Removal). *Let f be a maximum (S, T) -flow of G . If there exists two edges (v, t_1) and (v, t_2) with infinite capacity ($t_1, t_2 \in T$) we can either remove (v, t_1) or (v, t_2) from G without changing the amount of a maximum (S, T) -flow.*

Proof. Equivalent to proof of Lemma 5.2. \square

Definition 5.2 (Removable Edges). *We denote the set of all adjacent nodes v of a source node s resp. sink node t , where edge (s, v) or (v, t) is removable according to Lemma 5.2 and 5.3, with $\mathcal{R}(s)$ resp. $\mathcal{R}(t)$.*

The following lemma is a generalisation of Lemma 4.2. We will use the definition of $\text{in}(u)$ and $\text{out}(u)$ presented in Section 4.1. Further, $G_{V'}$ is a subgraph of $G = (V, E)$ induce by $V' \subseteq V$ (see Defintion 2.3).

Lemma 5.4 (General Source/Sink Node Removal). *Let f be a maximum (S, T) -flow of $G = (V, E, u)$ with $|f| < \infty$ and $E_s \subseteq \mathcal{R}(s)$ and $E_t \subseteq \mathcal{R}(t)$ with $s \in S$ and $t \in T$. If s is a source node where all outgoing edges have infinite capacity and t is a sink node where all incoming edges have infinite capacity, then $|f|$ is equal with the amount of a maximum (S', T) -flow of $G_{V \setminus \{s\}}$ and a maximum (S, T') -flow of $G_{V \setminus \{t\}}$, where $S' = (S \setminus \{s\}) \cup (\text{out}(s) \setminus E_s)$ and $T' = (T \setminus \{t\}) \cup (\text{in}(t) \setminus E_t)$.*

Proof. E_s is an arbitrary subset of $\mathcal{R}(s)$, where foreach $v \in E_s$ the edge (s, v) is removable. S' is the source set without node s extended with all outgoing edges of s minus the removable edges E_s . With Lemma 5.2 we can remove all edges (s, v) with $v \in E_s$ from G and obtain flow network G' . Finally, we can apply Lemma 4.2 on G' and obtain $G_{V \setminus \{s\}}$ with (S', T) as source and sink set (see Figure 18). All used Lemma's did not change the amount of a maximum flow. Therefore, a maximum (S, T) -flow of G is equal with a maximum (S', T) -flow of $G_{V \setminus \{s\}}$. The proof for t is equivalent. \square

The proof of Lemma 5.1 can be applied one-to-one on our new source and sink sets because $S_1 \subseteq S$ and $T_1 \subseteq T$. Therefore, S and T as defined in Equation 5.3 and 5.4 satisfies condition (i) of Problem 5.1. We will show that for S and T the equality $\Delta_H = \Delta_{H_{V'}}$ holds.

Theorem 5.1. *Let $B_1 = (V_1, V_2)$ be a bipartition of H and $T_L(H_{V'})$ the flow network of subhypergraph $H_{V'}$ with S and T as defined in Equation 5.3 and 5.4 (with $V' \subseteq V$). If B_2 is a bipartition obtained by a maximum (S, T) -flow computation on $T_L(H_{V'})$ with f as maximum flow, then $\omega_{H_{V'}}(B_2) = |f|$ ($\Rightarrow \Delta_H = \Delta_{H_{V'}}$).*

Proof. Consider the extension $Ex(H_{V'})$ of subhypergraph $H_{V'}$ (see Definition 5.1). Each maximum (S, T) -flow f' of $T_L(Ex(H_{V'}))$ is then equal with a minimum-weight (S, T) -cutset of $H_{V'}$ according to our definition of the cut $\omega_{H_{V'}}(B_2)$ over the cut hyperedges of H . Because each hyperedge which is partially contained in $H_{V'}$ is fully contained in $Ex(H_{V'})$. Therefore, it holds that $|f'| = \omega_{H_{V'}}(B_2)$. However, we have to model some restrictions into our source and sink set of $T_L(Ex(H_{V'}))$. We will denote the source and sink set of $T_L(Ex(H_{V'}))$ with S' and T' . Each hypernode contained in a non-cut border hyperedge $e \in \delta B_1$ should not be able to move such that we ensure that e is not cut after a maximum (S', T') -flow calculation. Therefore, we add S_1 and T_1 to S' and T' (see Equation 5.1 and 5.2). Further, all hypernodes $v \in V''$ (see Defintion 5.1) are not contained in $H_{V'}$. Consequently, they cannot change their block if we calculate a maximum (S, T) -flow of $T_L(H_{V'})$. Therefore, we add $V'' \cap V_1$ to S' and $V'' \cap V_2$ to T' . With S' and T' as source and sink set we ensure that only hypernodes $v \in V'$ are able to move and since $S_1 \subseteq S'$ and $T_1 \subseteq T'$, we ensure that $\omega_H(B_2) \leq \omega_H(B_1)$ (see Lemma 5.1).

In the following, we apply Lemma 5.4 on all hypernodes $v \in V''$ such that the flow network $T_L(Ex(H_{V'}))$ converge against $T_L(H_{V'})$ with S and T as source and sink set as defined in Equation 5.3 and 5.4 without changing the amount of the maximum flow f' . Since $|f'| = \omega_{H_{V'}}(B_2)$, for a maximum (S, T) -flow f of $T_L(H_{V'})$ then holds $|f| = \omega_{H_{V'}}(B_2)$. Per definition a node $v \in V''$ is either a source or sink node. For each source node $s \in V''$ we have to define a removable subset $E_s \subseteq \mathcal{R}(s)$ such that after removing all $s \in S' \cap V''$ the resulting source set

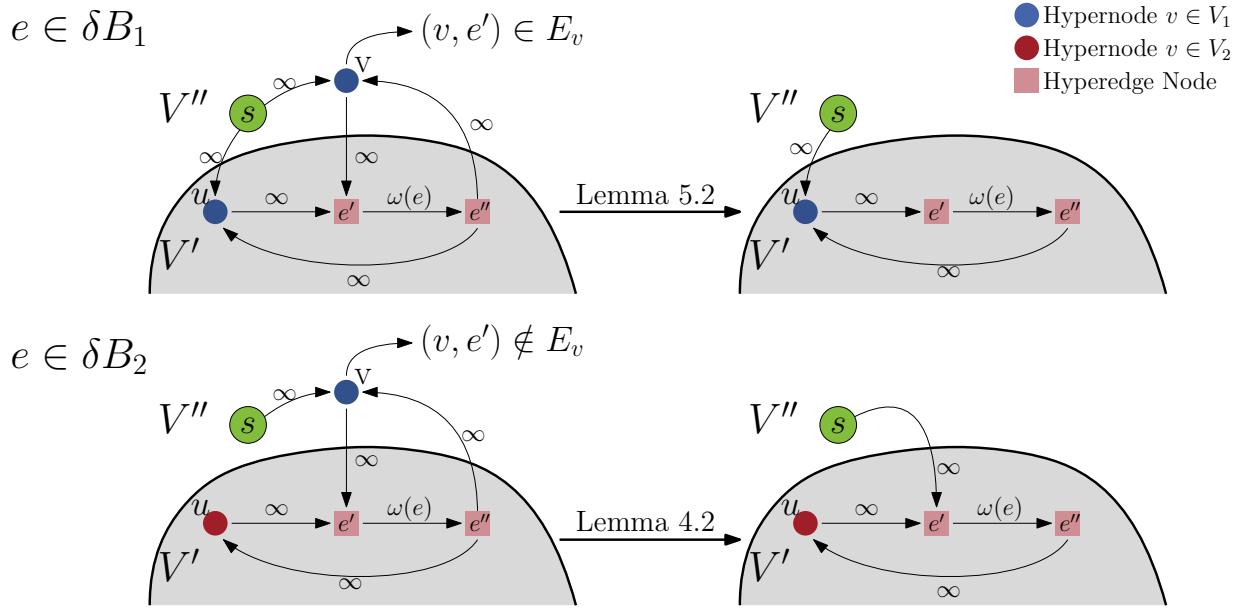


Figure 19: Illustration how to remove a source node $v \in V''$. Note, the green node s is the super source of the flow problem. Consequently, all nodes connected to s are source nodes in the corresponding *multi-source multi-sink* flow problem.

S is equal to Equation 5.3. The technique for removing each sink node $t \in S' \cap V''$ will be equivalent. For a source node $s \in V''$ we define $E_s = \{e' \mid e \in I(s) \cap \delta B_1\}$. Remember, δB_1 contains all non-cut border hyperedges of H . Thus, for each $e \in I(s) \cap \delta B_1$ exists a source node $\bar{s} \in V'$ of S_1 such that edges (s, e') and (\bar{s}, e') are contained in $T_L(Ex(H_{V'}))$ (see Figure 19). Therefore, E_s is a removable subset of $\mathcal{R}(s)$. For each $t \in T' \cap V''$ we define the removable subset $E_t = \{e'' \mid e \in I(t) \cap \delta B_1\}$. Applying Lemma 5.4 on all source nodes $s \in V''$ with E_s as removable subset and on all sink nodes $t \in V''$ with E_t as removable subset yield flow network $T_L(H_{V'})$ with S and T as source and sink set as defined in Equation 5.3 and 5.4. A hyperedge $e \in \delta B_1$ cannot become a source or sink node because if we remove a source node $s \in e$ then e' is in the removable subset E_s (see Figure 19). The same holds for each sink node $t \in V''$. A hyperedge $e \in \delta B_2$ becomes a source node of S if we remove a source node $s \in e$ and a sink node of T if we remove a sink node $t \in e$ because $\forall e \in \delta B_2 : e', e'' \notin E_s \cup E_t$. Therefore, S and T are equal to our source and sink set definition and for a maximum (S, T) -flow f it holds that $|f| = |f'| = \omega_{H_{V'}}(B_2)$. \square

We are now able to extract a subhypergraph $H_{V'}$ out of an already bipartitioned hypergraph H and calculate a minimum (S, T) -bipartition of $H_{V'}$ with S and T as defined in Equation 5.3 and 5.4. The resulting bipartition induce a new cut on H smaller or equal than the old cut. Further, we show with our modeling technique of S and T that Δ_H can be calculated with the help of the amount of a maximum (S, T) -flow computation on $T_L(H_{V'})$.

In Section 4.3 we described how to remove hyperedges of size $|e| = 2$ by adding an undirected flow edge between the corresponding vertices $u, v \in e$. However, if the incoming or outgoing hyperedge node is a source or a sink node, we can not directly remove the hyperedge nodes. There are two special cases which are illustrated in Figure 21. This situation occurs if one of the two vertices is part of the flow problem and one not. In case, if the incoming hyperedge node e' is a source node, we only remove the outgoing hyperedge node e'' and add a directed flow edge from e' to v with capacity $\omega(e)$. In the second case, if the outgoing hyperedge node e'' is a sink node, we only remove the incoming hyperedge node e' and add a directed flow edge from v to e'' with capacity $\omega(e)$.

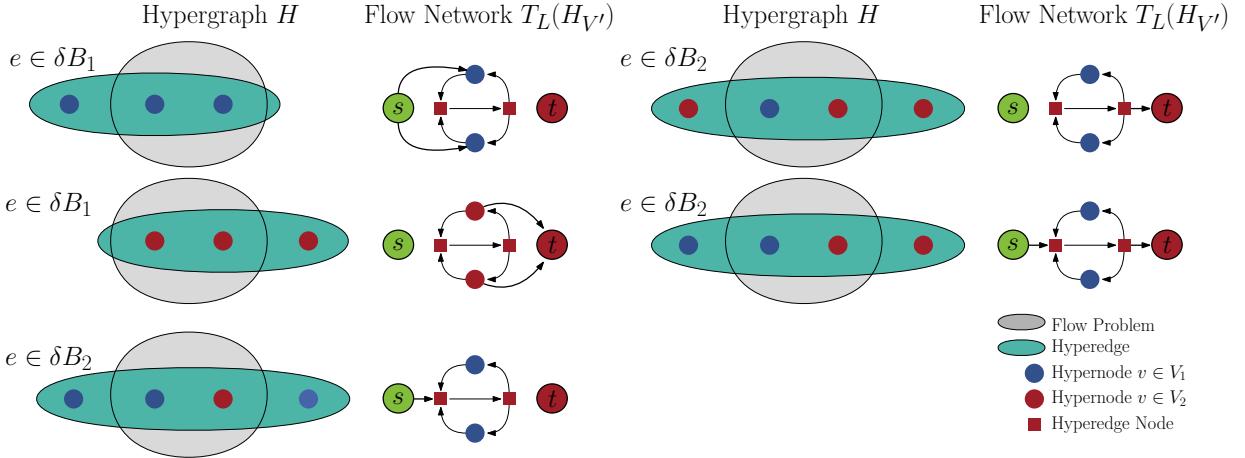


Figure 20: Illustration of modeling sources and sinks defined in Equation 5.3 and 5.4.

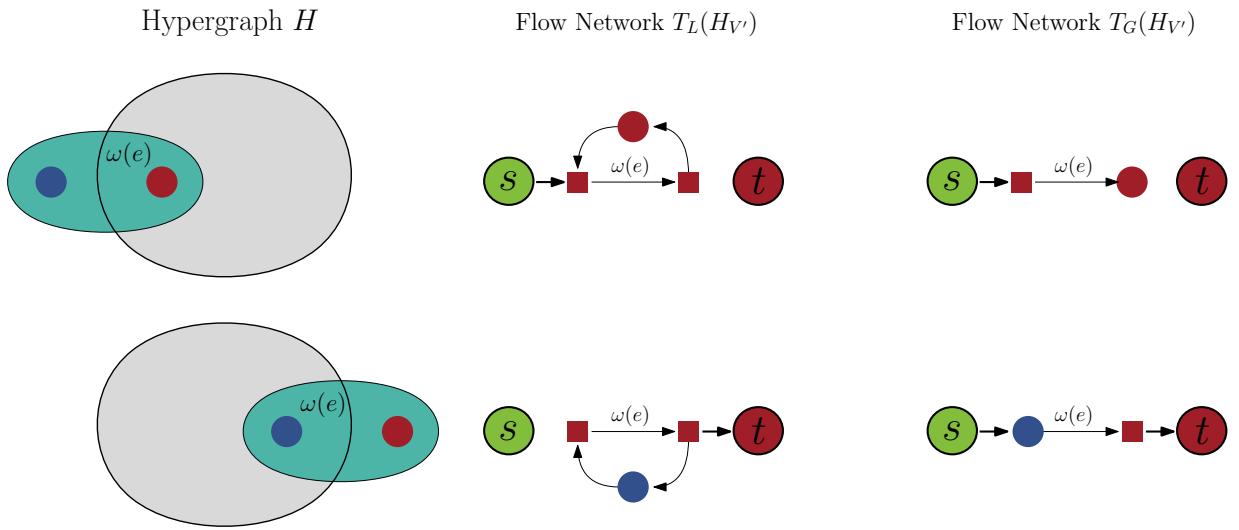


Figure 21: Illustration of modeling hyperedges of size two if the incoming or outgoing hyperedge node is a source or a sink node of the flow problem.

With the given approach we can optimize the cut metric of a given bipartition of a hypergraph H . We can transfer those results to improve a k -way partition $\Pi = (V_1, \dots, V_k)$ if the objective is the connectivity metric. Let $V' \subseteq V_i \cup V_j$ be a subset of the hypernodes of two adjacent blocks V_i and V_j . If we optimize the cut of subhypergraph $H_{V'}$ we simultaneously optimize the connectivity metric of H . The reduction of the cut of $H_{V'}$ is then equal with the decrease in the connectivity metric of H .

5.2. Most Balanced Minimum Cuts on Hypergraphs

Picard and Queyranne [39] show that all minimum (s, t) -cuts of a graph G are computable with one maximum (s, t) -flow computation by iterating through all *closed node sets* of the residual graph of G . The corresponding algorithm is presented in Section 3.3.3.

We can apply the same algorithm on hypergraphs. A minimum-capacity (s, t) -cutset of $T_L(H)$ is equal with a minimum-weight (s, t) -cutset of H . With the algorithm of Section 3.3.3 we can find all minimum-capacities (s, t) -cutsets of $T_L(H)$, which are also minimum-weight (s, t) -cutsets of H . The corresponding minimum-weight (s, t) -bipartitions are all *closed node sets* of the residual graph of $T_L(H)$.

However, when we use e.g. $T_H(H, V')$ (see Section 4.1) or $T_{\text{Hybrid}}(H, V')$ (see Section 4.4) as underlying flow network, some hypernodes are removed from the flow problem. It is a problem if we want to enumerate all minimum-weight (s, t) -bipartitions. The solution for this problem is quite simple. After a maximum (s, t) -flow calculation on one of the two mentioned networks we insert all removed hypernodes with their corresponding edges again into the residual graph of our flow network. The maximum (s, t) -flow is still maximal. Otherwise, we would have found an *augmenting path* on the flow network before. We are now able to compute all minimum-weight (s, t) -bipartitions the same way as with $T_L(H)$.

5.3. A direct k -way Flow-Based Refinement Framework

We have described how a hypergraph H could be transformed into a flow network $T_L(H)$ such that each minimum-capacity (S, T) -cutset of $T_L(H)$ is a minimum-weight (S, T) -cutset of H (see Section 3.2). Additionally, we present techniques to sparsify the flow network $T_L(H)$ [30] to reduce the complexity of the flow problem (see Section 4). Further, we show how to configure the source and sink sets of a flow network of a subhypergraph $H_{V'}$ (with $V' \subseteq V$) such that a *Max-Flow-Min-Cut* computation improves a given bipartition of H (see Section 5.1). Finally, we can enumerate all minimum-weight (s, t) -cutsets of a subhypergraph $H_{V'}$ with one maximum (S, T) -flow calculation [39].

We will now present our direct k -way flow-based refinement framework which we integrated into the n -level hypergraph partitioner *KaHyPar* [22] (see Section 3.4.2). Our flow-based refinement approach optimizes the *connectivity* metric. We used a similar architecture as proposed by Sanders and Schulz [42] (see Section 3.3). The basic concepts of the framework are illustrated in Figure 22.

Our maximum flow calculations are embedded into an *Active Block Scheduling* refinement [23] (see Section 3.3.4). Each time we use flows to improve the connectivity metric of a given k -way partition Π we construct the quotient graph Q of Π . Afterwards, we iterate over all edges of Q in random order. For each edge (V_i, V_j) of Q , we build a flow problem around the cut of the bipartition induced by V_i and V_j . To do that we use two *BFS*, one only touches hypernodes of V_i and the second only touches hypernodes of V_j . The *BFS* is initialized with all hypernodes contained in a cut hyperedge of the bipartition (V_i, V_j) . A pairwise flow-based refinement is embedded into the *adaptive flow iterations* strategy [42] (see Section 3.3.2) which also determines the size of the flow problem.

After we define the subhypergraph $H_{V'}$, which we use to improve the bipartition (V_i, V_j) on H , we construct one of the flow networks proposed in Section 4 with sources S and sinks T defined in Section 5.1. We implemented two maximum flow algorithms. One is a slightly modified *augmenting path* algorithm of Edmond & Karp [15] (see Section 3.1.1) and the second is the *Push-Relabel* algorithm of Goldberg & Tarjan [11, 20] (see Section 3.1.2). Since we have a *Multi-Source-Multi-Sink* problem, we can find several *augmenting paths* with one *BFS*. After we execute a *BFS* on the residual graph, we search as many as possible edge-disjoint paths in the resulting *BFS*-tree connecting a source s with a sink t . Our Goldberg & Tarjan implementation uses a *FIFO* queue and the *global relabeling* and *gap* heuristic [11]. We do not use an external implementation of a maximum flow algorithm. Since the $I|O$ of writing a flow problem to memory and reading the solution would significantly slowdown the performance of our algorithm because we have to solve an enormous number of flow problems during the *Active Block Scheduling* refinement. After we determine a maximum (S, T) -flow on our flow network, we iterate over all minimum (S, T) -bipartitions of $H_{V'}$ [39] and choose the *Most Balanced Minimum Cut* (see Section 3.3.3 and 5.2) according to our *balanced constraint*.

KaHyPar is an n -level hypergraph partitioner ($|V| = n$) taking the multilevel paradigm to its

extreme by removing only a single vertex in every level of the hierarchy [1] (see Section 3.4.2). During the refinement step n local searches are instantiated. Therefore, using our flow-based refinement as local search algorithm on each level is not applicable, because the performance slowdown would be tremendous. Therefore, we introduce *Flow Execution Policies*. One is to execute our flow-based refinement on each level i where $i = \beta \cdot j$ with $j \in \mathbb{N}_+$ and β as a predefined tuning parameter. Another approach is to simulate a multilevel partitioner with $\log(n)$ hierarchies. A flow-based refinement is then executed on each level i where $i = 2^j$ with $j \in \mathbb{N}_+$. Each policy also performs the *Active Block Scheduling* refinement on the last level of the hierarchy. In all remaining levels where no flow is executed, we can use an *FM*-based local search algorithm [1, 16, 41] (see Section 3.3.4).

An observation during the implementation of this framework was that only a minority of the pairwise refinements based on flows yields to an improvement of the connectivity metric on hypergraph H . Thus, we introduce several rules which might prevent unnecessary flow executions to improve the effectiveness ratio by simultaneously speeding up the running time.

- (R1) If a flow-based refinement did not lead to an improvement on two blocks in all previous executions, we would use flows only in the first iteration of *Active Block Scheduling*.
- (R2) If the cut between two adjacent blocks in the quotient graph is small (e.g. ≤ 10) we skip the flow-based refinement on these blocks except on the last level of the hierarchy.
- (R3) If the value of the cut of a minimum (S, T) -bipartition of $H_{V'}$ is the same as the cut before, we stop the pairwise refinement.

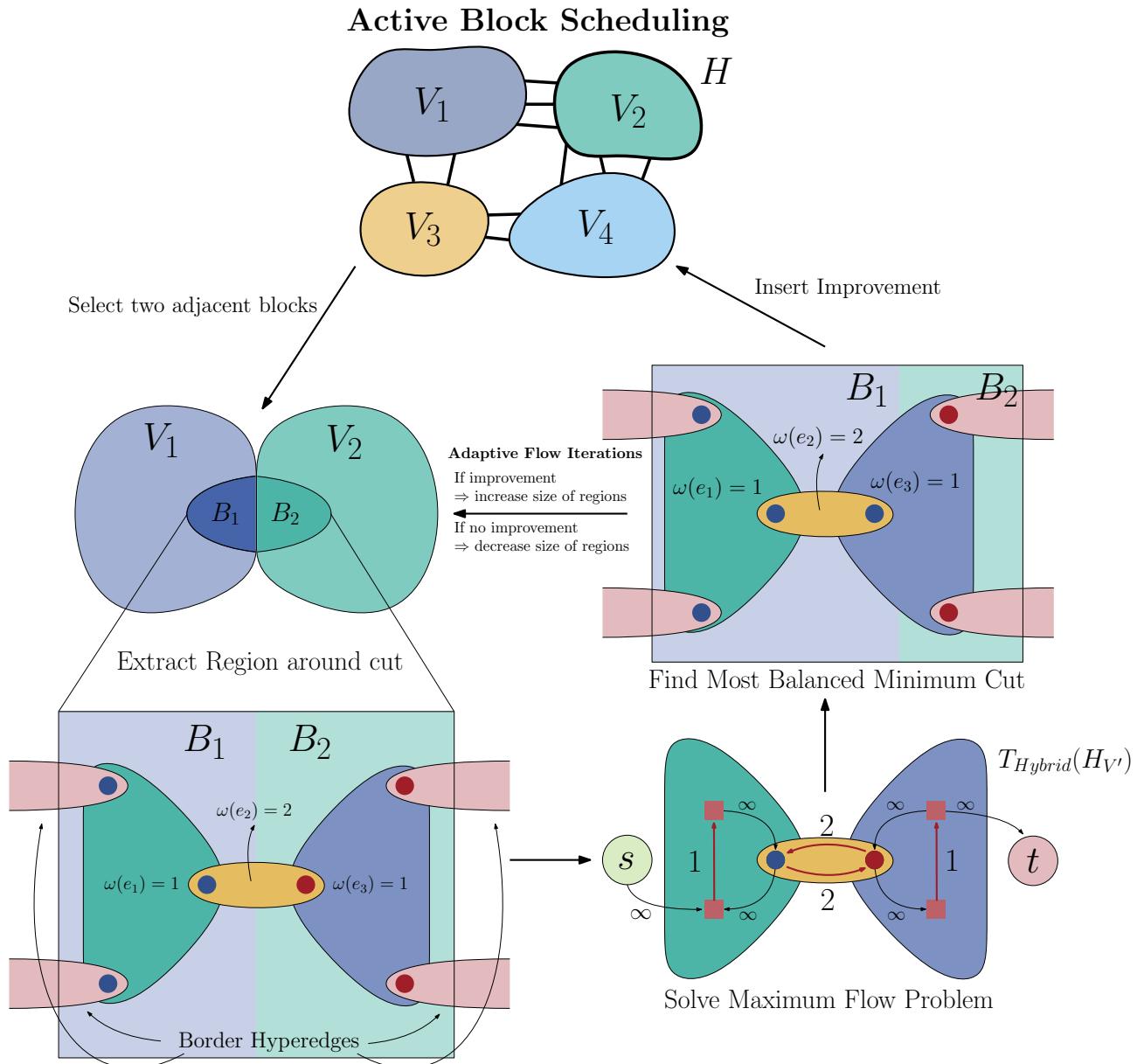


Figure 22: Illustration of our flow-based refinement framework for direct k -way hypergraph partitioning.

6. Experimental Results

In this Section, we evaluate the performance of our flow-based refinement framework proposed in Section 4 and 5. We examine the impact of our sparsifying techniques of the *Lawler-Network* [30] on the performance of a maximum flow algorithm (see Section 6.3). Further, several configurations with different heuristics enabled or disabled are compared against the baseline configuration of *KaHyPar* to optimally configure our flow-based refinement algorithm (see Section 6.4 and 6.5). Finally, we compare our final configuration against other state-of-the-art hypergraph partitioners (see Section 6.6).

6.1. Instances

Our full benchmark set consists of 488 hypergraphs. We choose our benchmarks from three different research areas. For VLSI design we use instances from the *ISPD98 VLSI Circuit Benchmark Suite* (ISPD98) [2] and add more recent instances of the *DAC 2012 Routability-Driven Placement Contest* (DAC) [46]. Further, we interpret the Sparse Matrix instances of the *Florida Sparse Matrix collection* (SPM) [12] as hypergraphs using the row-net model [9]. The rows of each matrix are treated as hyperedges and the columns are the vertices of the hypergraph. Our last benchmark type are SAT formulas of the *International SAT Competition 2014* [6]. A common interpretation of a SAT formula as hypergraph is to interpret the literals as vertices and each clause as a net (LITERAL) [37]. Mann and Papp [33] suggested two other hypergraph representation of SAT formulas, called PRIMAL and DUAL. The PRIMAL representation treats each variable as vertex and each clause as hyperedge. The DUAL representation treats each clause as vertex and the variables induced nets containing all clauses where the corresponding variable occurs. A statistical summary of the different instance types is presented in Table 7.

We divide our full benchmark set into two smaller subsets. Our *parameter tuning* benchmark set consists of 25 hypergraphs, 5 of each instance type (except DAC). Additionally, we choose a benchmark subset of 165 instances. On our general experiments we partition each hypergraph into $k \in \{2, 4, 8, 16, 32, 64, 128\}$ blocks and use for each k 10 different *seeds* with $\epsilon = 3\%$.

6.2. System and Methodology

Our experiments run on a single core of a machine consisting of two *Intel Xeon E5- 2670 Octa-Core* processors clocked at 2.6 GHz. The machine has 64 GB main memory, 20 MB L3- and 8×256 KB L2-Cache. The code is written in C++ and compiled using g++-5.2 with flags `-O3 -mtune=native -march=native`. We refer to our new implementation of *KaHyPar* with (*M*)ax-(*F*)low-Min-Cut computations as *KaHyPar-MF* and the latest configuration with (*C*)ommunity-(*A*)ware coarsening as *KaHyPar-CA*.

We compare *KaHyPar-MF* against the state-of-the-art hypergraph partitioner *hMetis* [26, 27] and *PaToH* [9]. *hMetis* provides a direct k -way (*hMetis-K*) and recursive bisection (*hMetis-R*) implementation. Further, we also use the default configuration (*PaToH-D*) and quality preset (*PaToH-Q*) of *PaToH*. We configure *hMetis* to optimize the *sum-of-external-degree-metric* (SOED) and calculate $(\lambda - 1)(\Pi) = \text{SOED}(\Pi) - \text{cut}(\Pi)$. This is also suggested by the authors of *hMetis* [27]. Further, we have to adapt the imbalance definition of *hMetis-R*. An imbalance value of 5 means that the weight of each bisected block is allowed to be between $0.45 \cdot c(V)$ and $0.55 \cdot c(V)$. To ensure that *hMetis-R* produces a valid ϵ -balanced partition after $\log_2(k)$

bisections we have to adapt ϵ to

$$\epsilon' = 100 \cdot \left(\left((1 + \epsilon) \frac{\lceil \frac{c(V)}{k} \rceil}{c(V)} \right)^{\frac{1}{\log_2(k)}} - 0.5 \right)$$

If we evaluate the performance of our hypergraph partitioner, we first calculate the average (or minimum) of the different *seeds* of a hypergraph instance and then the *geometric mean* between all instances to give every instance comparable influence on the final result. To compare the performance of different hypergraph partitioner more detailed we use performance plots introduced in [43]. For each partitioner P and instance H we calculate the values $q_{H,P} := 1 - \text{best}_H/\text{algorithm}_{H,P}$ where best_H is the best quality achieved by a partitioner for instance H and $\text{algorithm}_{H,P}$ refers to the quality achieved by partitioner P for instance H . Afterwards, we sort all values $q_{H,P}$ of a partitioner P in decreasing order. For each partitioner P we plot the points $(H, q_{H,P})$. The faster the $q_{H,P}$ values intersect the zero line the better the performance of a partitioner in comparison to the others. If a partition of a partitioner P is not ϵ -balanced we set $q_{H,P} = 1 + \beta$ (with $\beta > 0$).

6.3. Flow Algorithms and Networks

In the first experiment, we want to examine the impact of our sparsifying techniques (see Section 4) on the performance of our maximum flow algorithms GOLDBERG-TARJAN and EDMOND-KARP. Therefore, we first take a look at the reduction of the number of nodes and edges on different benchmark types when using T_L (see Section 3.2), T_H (see Section 4.2), T_G (see Section 4.3) and T_{Hybrid} (see Section 4.4). Further, we want to evaluate the performance of the two implemented maximum flow algorithms on these networks.

We evaluate the performance of the different flow networks on flow problems with size $|V'| \in \{500, 1000, 5000, 10000, 25000\}$ hypernodes. The instances are generated by executing *KaHyPar* on our benchmark subset (see Table 6) for $k = 2$ and five different seeds. After an instance is bipartitioned, we generate flow problem instances with the above-mentioned sizes and execute each possible combination of flow algorithm and network on it.

The benchmark instances can be split into 6 different benchmark types. The properties of these instances regarding the average hypernode degree and average hyperedge size is shown in Table 6. Remember, T_G should perform best on instances with a small average hyperedge size and T_H should perform best on instances with a low average hypernode degree. Based on Table 6, T_G should significantly reduce the number of nodes and edges on PRIMAL and LITERAL instances and T_H on DUAL instances in comparison to our baseline T_L . Also both should sparsify the resulting flow network of ISPD98 and DAC instances. Further, we expect that T_{Hybrid} combines the advantages of both networks and performs best on all benchmark instances.

Figure 23 shows the predicted behavior for flow problems of size 25000 hypernodes. T_{Hybrid} reduces the number of nodes of nearly every benchmark type by at least a factor of 2, except on SPM instances. Another observation is that instances with a large average hypernode degree, like PRIMAL or LITERAL, yield to big flow problem instances and vice versa (see DUAL instances).

In Figure 24 we compare the performance of our flow algorithms on different flow networks. The bars in the plot indicates speedups relative to the flow algorithm EDMOND-KARP on flow network T_L . The main observation is that EDMOND-KARP performs better on small flow network instances and GOLDBERG-TARJAN on large flow network instances. For $|V'| \leq 1000$ EDMOND-KARP is faster than GOLDBERG-TARJAN in most of the different benchmark types.

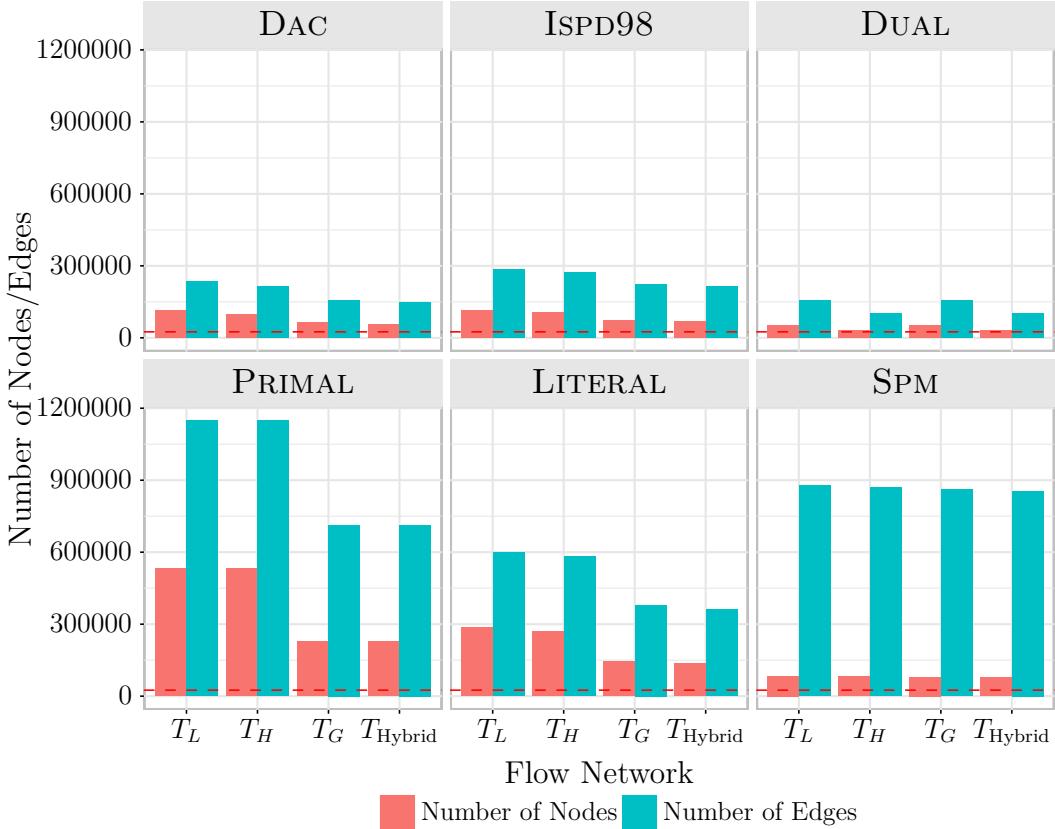


Figure 23: Comparison of the number of nodes and edges induced by flow problems of size $|V'| = 25000$ on our flow network for different benchmark types. The red dashed lines indicates 25000 nodes.

For $|V'| > 1000$ we can observe the opposite behavior except for **DAC** and **DUAL** instances. But the resulting flow problems of these instances are still the smallest among all benchmark types (see Figure 23). On the largest flow network instances **PRIMAL** and **LITERAL** for $|V'| = 25000$ GOLDBERG-TARJAN is up to a factor of 4-7 faster than EDMOND-KARP. Further, both algorithms perform best on T_{Hybrid} . Table 1 shows the summary of our flow algorithm and network experiment on all benchmark instances. It proofs our assumption that EDMOND-KARP works best on small instances and GOLDBERG-TARJAN on large instances. However, our *Max-Flow-Min-Cut* computations are embedded in an *Adaptive Flow Iteration* strategy (see Section 3.3.2). Therefore, the running time of flow instances generated with a large α will dominate the ones with small α . Thus, we choose GOLDBERG-TARJAN in combination with our flow network T_{Hybrid} in the following experiments.

6.4. Setup of the direct k -way Flow-Based Refinement

In this Section, we examine the quality of our k -way flow-based refinement algorithm with different configurations on our parameter tuning benchmark subset (see Table 5). There are several configurations and tuning parameters which we have to evaluate:

- *Max-(F)low-Min-Cut* computations as refinement algorithm (see Section 5.3)
- *Adaptive Flow Iteration* parameter α' (see Section 3.3.2)
- *(M)ost Balanced Minimum Cut* heuristic (see Section 5.2)
- Combining *Max-(F)low-Min-Cut* computations with *(FM)* refinement

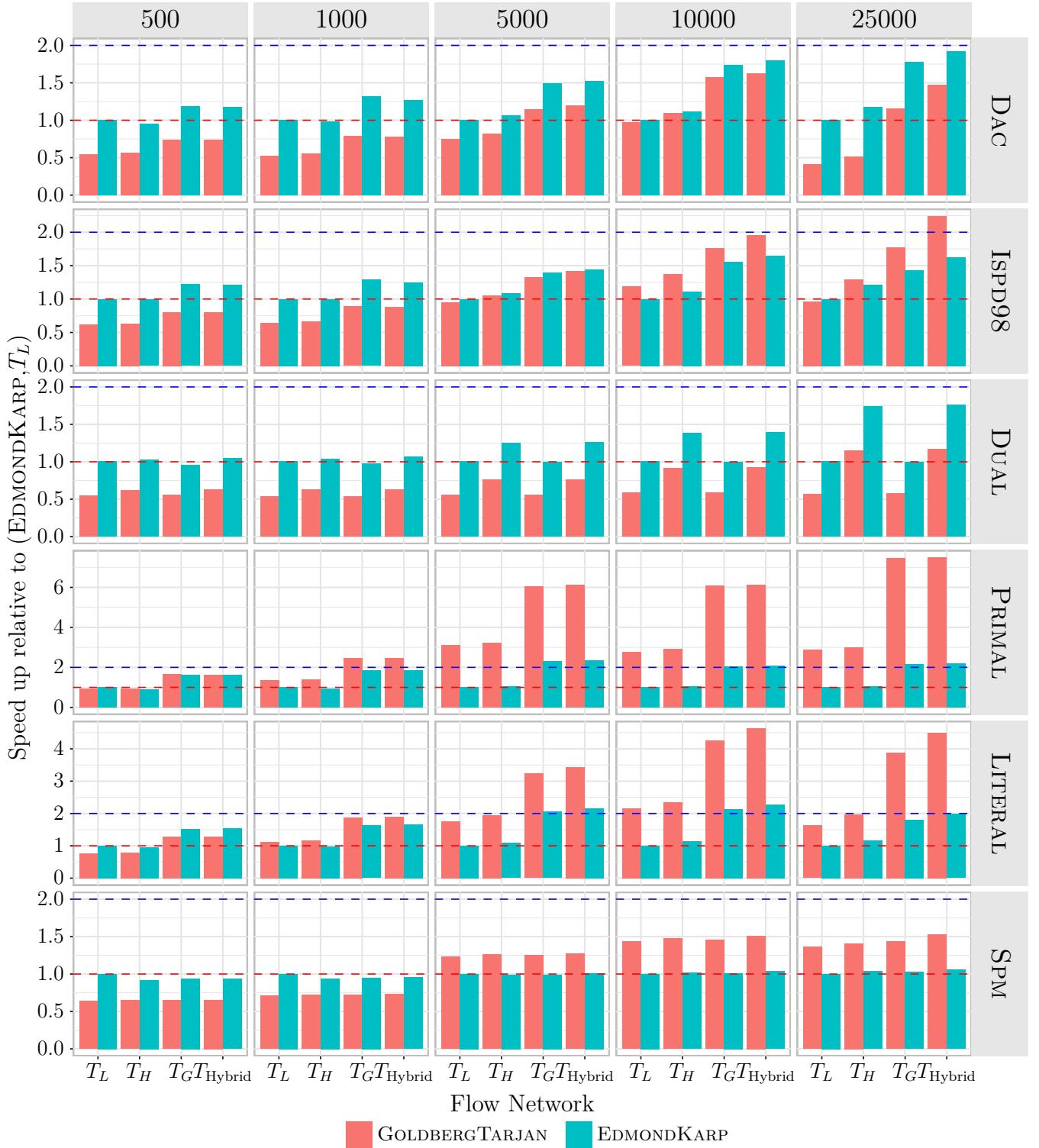


Figure 24: Speedup of our flow algorithms and networks relative to EDMONDKARP on T_L for different instance sizes and types. The red dashed line indicates the (EDMONDKARP, T_L) implementation and the blue dashed line indicates a speedup by a factor of 2.

Instance	GOLDBERG-TARJAN				EDMOND-KARP			
	$ V' $	T_{Hybrid}	T_G	T_H	T_L	T_{Hybrid}	T_G	T_H
		$t[ms]$	$t[\%]$	$t[\%]$	$t[\%]$	$t[\%]$	$t[\%]$	$t[\%]$
ALL	500	0.91	+2.24	+24.93	+29.35	-25.39	-24.3	-6.68
	1000	1.95	+3.65	+26.19	+32.95	-13.99	-12.36	+10.81
	5000	13.71	+8.63	+29.39	+43.11	+27.03	+35.33	+73.97
	10000	30.54	+12.57	+36.15	+54.62	+47.93	+61.72	+100.41
	25000	67.96	+23.36	+52.12	+87.8	+53.25	+77.85	+100.95

Table 1: Running time comparison of maximum flow algorithms on different flow networks.

Note, all values in the table are in percentage relative to GOLDBERG-TARJAN on flow network T_{Hybrid} . In each line the fastest variant is marked bold.

In the following, we will denote a configuration e.g. with (+F,-M,-FM) which indicates which heuristic resp. technique is enabled (+) or disabled (-). The meaning of the abbreviations is explained in the enumeration above (see letters inside parenthesis). We evaluate a configuration for $k \in \{2, 4, 8, 16, 32, 64, 128\}$, $\alpha' \in \{1, 2, 4, 8, 16\}$ and 10 different seeds on our parameter tuning benchmark subset ($\epsilon = 3\%$). Our pairwise flow-based refinement is embedded in a k -way *Active Block Scheduling* refinement which is executed on each level i with $i = 2^j$ ($j \in \mathbb{N}_+$) (see Section 5.3). Additionally, we tested configuration (+F,+M,+FM) with *flow execution policy* $i = 128j$. This configuration has an impracticable running time, but should provide a lower bound. We refer to this variant as CONSTANT128. As a baseline reference, we use the latest quality configuration of *KaHyPar* (KaHyPar-CA) [22].

The results are summarized in Table 2. The values in the column *Avg* are improvements of the connectivity metric relative to our baseline configuration (-F,-M,+FM). The running time are absolute values in seconds. The first observation is that flows on its own as refinement strategy are not strong enough to outperform the *FM* heuristic. Our strongest configuration with $\alpha' = 16$ is 2.58% worse than our *FM* baseline. But the result is still remarkable because we only execute flows on $\log n$ levels instead of n as the *FM* algorithm does. The running time scales nearly linear with parameter α' .

Enabling the *Most Balanced Minimum Cut* heuristic significantly improves the quality compared to the baseline flow configuration (+F,-M,-FM). But the quality improvements are more significant for large α' . The larger the flow problem, the larger is the number of different minimum (S, T) -cutsets and this increases the possibility to find a feasible solution according to our balanced constraint. Also it outperforms our baseline *FM* configuration for $\alpha' = 16$ by 0.51%. If we enable *FM* refinement at all levels where no flow is executed, we improve the solution quality by nearly 2% (for $\alpha' = 16$). Also, the running time of this variant is faster than all previous flow configurations because we transfer more work to the *FM* refinement. It has as a consequence that a block becomes faster *inactive* during *Active Block Scheduling* and this decreases the number of rounds of complete pairwise flow-based refinements on the quotient graph.

Finally, CONSTANT128 gives us a lower bound of the quality achievable with a combination of flow-based and *FM* refinement. Flows are executed in each 128th level of the multilevel hierarchy. The quality is 2.44% better than our baseline configuration, but ≈ 100 slower. Compared to (+F,+M,+FM) for $\alpha = 16$, CONSTANT128 is only 0.57% better and around ≈ 25 times slower.

Our best configuration is (+F,+M,+FM) with $\alpha' = 16$. It is also the most effective one

Config.	(+F,-M,-FM)		(+F,+M,-FM)		(+F,+M,+FM)		CONSTANT128	
α'	Avg.[%]	$t[s]$	Avg.[%]	$t[s]$	Avg.[%]	$t[s]$	Avg.[%]	$t[s]$
1	-15.48	12.94	-15.26	13.29	0.14	14.99	0.32	67.38
2	-10.5	16.07	-10.12	16.93	0.36	16.93	0.62	139.21
4	-5.98	21.22	-5.08	23.01	0.67	20.76	1.03	274.6
8	-3.22	30.73	-1.64	33.72	1.25	28.65	1.67	558.81
16	-1.52	50.89	0.51	56.39	1.87	46.17	2.44	1220.92
Ref.	(-F,-M,+FM)		6373.88	13.73				

Table 2: Table contains results for different configurations of our flow-based refinement framework for increasing α' . The quality in column *Avg.* is relative to our baseline configuration without the usage of flows.

(see Effectiveness Test in Appendix C). For further experiments, we refer to this variant as KaHyPar-MF.

6.5. Speed-Up Heuristics

At the end of Section 5.3, we present several heuristics to prevent unnecessary flow executions during *Active Block Scheduling* ((R1)-(R3)). The main assumption is that only a minority of *Max-Flow-Min-Cut* computations lead to an improvement on H . To prove that we execute KaHyPar-MF on our benchmark subset (see Table 6) and enable one heuristic after another. Table 3 summarizes the results of the experiment. KaHyPar-CA is the currently best configuration of *KaHyPar* and KaHyPar-MF is our baseline flow configuration of Section 6.4. The index of the remaining variants of KaHyPar-MF describes which speed-up heuristics are enabled (see Section 5.3). On average, enabling all speed-up heuristics worsen the quality of KaHyPar-MF only by 0.09%. On the other hand, the *Max-Flow-Min-Cut* computations are significantly faster by a factor of ≈ 2 . In its final configuration KaHyPar-MF_(R1,R2,R3) computes partitions with 2% better quality ($(\lambda - 1)$ -metric) than KaHyPar-CA by a slowdown only of a factor of ≤ 2 . In the following, we will denote our final configuration KaHyPar-MF_(R1,R2,R3) with KaHyPar-MF.

Variant	Avg.[%]	Min.[%]	$t_{\text{flow}}[s]$	$t[s]$
KaHyPar-CA	7077.2	6820.17	-	29.26
KaHyPar-MF	-2.13	-1.8	52.28	81.54
KaHyPar-MF _(R1)	-2.05	-1.74	41.48	70.74
KaHyPar-MF _(R1,R2)	-2.05	-1.73	35.27	64.54
KaHyPar-MF _(R1,R2,R3)	-2.04	-1.75	27.62	56.88

Table 3: Results of our flow-based refinement framework with different speedup heuristics.

6.6. Comparison with other Hypergraph Partitioner

Finally, we compare our new approach KaHyPar-MF with different state-of-the-art hypergraph partitioner on our full benchmark set. We excluded 194 instances of 3416 either because PaToH-

Q could not allocate enough memory or other partitioners did not finish in time. The excluded instances are shown in Table 9.

Figure 25 summarizes the results of the experiment. KaHyPar-MF produced on $\approx 70\%$ of all benchmark instances the best partition. It is followed by hMetis-R (14%), hMetis-K (11%), KaHyPar-CA (2.4%), PaToH-Q (1.9%) and PaToH-D (1.4%). Since KaHyPar-MF builds on top of KaHyPar-CA, it outperforms KaHyPar-CA on most of the instances. Comparing KaHyPar-MF individually with each partitioner, KaHyPar-MF produced better partitions than KaHyPar-CA, hMetis-R, hMetis-K, PaToH-Q, PaToH-Q in 96%, 80%, 82%, 95%, 95% cases. Especially on *VLSI* instances, KaHyPar-MF calculates significantly better partitions than all other hypergraph partitioners (see DAC and ISPD98 in Figure 25).

Table 12 shows the running time of all partitioner on different benchmark types. The running time of KaHyPar-MF is within a factor of 2 slower than KaHyPar-CA and is comparable to the running time of hMetis-K.

Partitioner	Running Time $t[s]$						
	ALL	DAC	ISPD98	PRIMAL	LITERAL	DUAL	SPM
KaHyPar-MF	62.24	637.58	22.29	71.63	140.84	106.24	29.61
KaHyPar-CA	31.05	368.97	12.35	32.91	64.65	68.27	13.91
hMetis-R	79.23	446.36	29.03	66.25	142.12	200.36	41.79
hMetis-K	57.86	240.92	23.18	44.23	94.89	125.55	35.95
PaToH-Q	5.89	28.34	1.89	6.9	9.24	10.57	3.42
PaToH-D	1.22	6.45	0.35	1.12	1.58	2.87	0.77

Table 4: Comparing the average running time of KaHyPar-MF with KaHyPar-CA and other hypergraph partitioners.

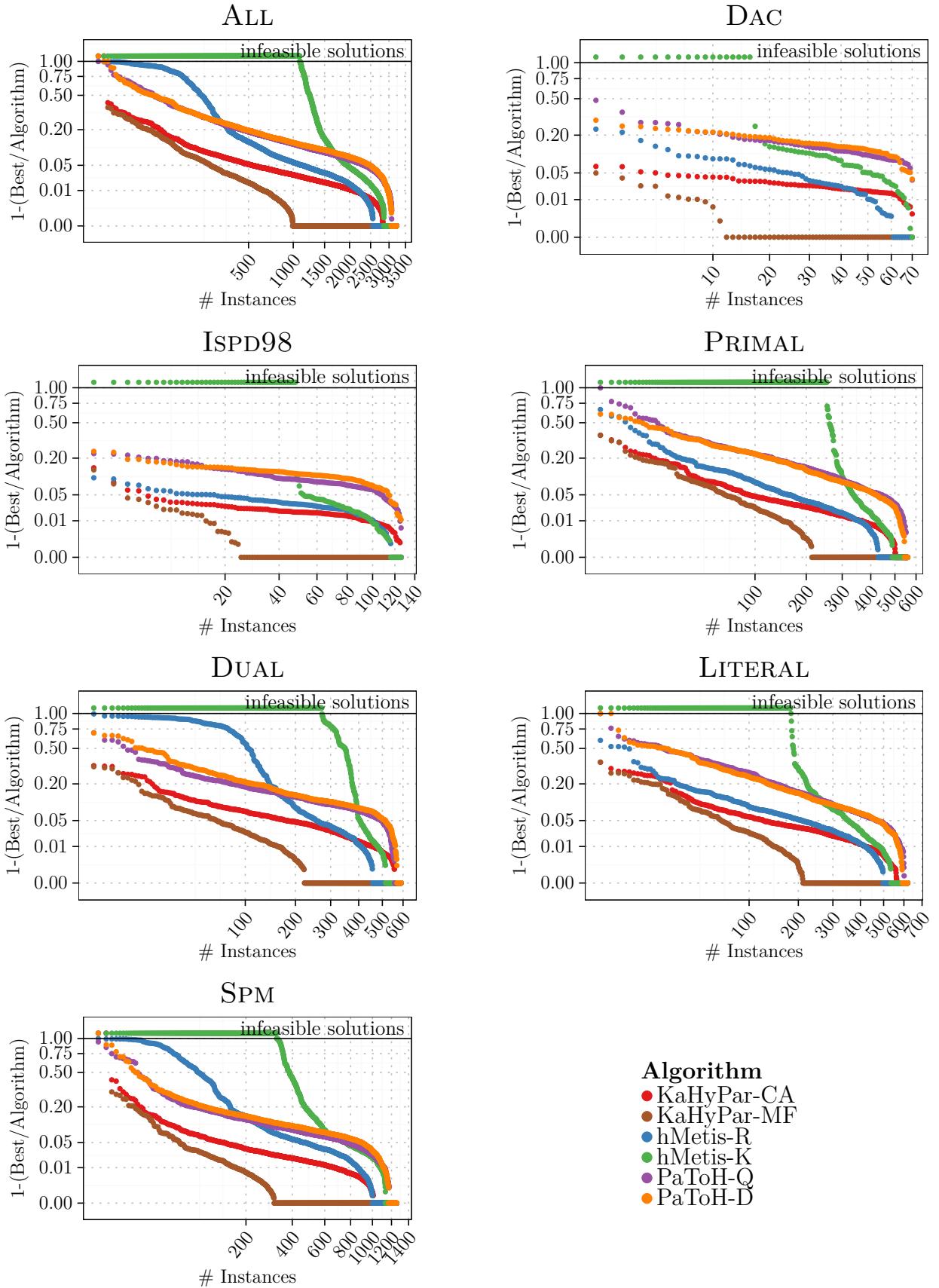


Figure 25: Min-Cut performance plots comparing KaHyPar-MF with KaHyPar-CA and other systems. Plots are explained in Section 6.2.

7. Conclusion

In this thesis, we developed a novel *local search* technique based on *Max-Flow-Min-Cut* computations for multilevel hypergraph partitioning. We integrated our *flow-based refinement* framework into the n -level hypergraph partitioner *KaHyPar* and show that in combination with the *FM* heuristic our new approach produces the best-known partitions for a wide range of applications.

On the road to a practical implementation, we developed several concepts to speed up flow computations on a flow network of a hypergraph (see Section 3.2). One is to remove low-degree hypernodes from the network and instead insert a clique between all incident hyperedge nodes. We show that the number of nodes and edges could be reduced if the degree of a hypernode is smaller or equal than 3. Further, we model a hyperedge of size 2 as an undirected flow edge. We combine both techniques in a *Hybrid-Network* and show that maximum flow algorithms are up to a factor of 3 faster compared to the execution on the *Lawler-Network* [30] on real-world benchmarks.

Our *flow-based refinement* framework is based on the ideas of Sanders and Schulz [42] (developed for multilevel graph partitioning). Given an already bipartitioned hypergraph, we show how to configure the source and sink sets of the flow network of a subhypergraph such that a *Max-Flow-Min-Cut* computation yields a cut smaller or equal than the cut before on the original hypergraph. Further, we extended the source and sink sets with additional nodes such that we can calculate the cut of hypergraph H after a *Max-Flow-Min-Cut* computation with the help of the amount of a maximum (S, T) -flow on a subhypergraph $H_{V'}$ in constant time. Additionally, we explain how one can find all minimum (s, t) -cutsets with one maximum (s, t) -flow calculation on hypergraphs.

We integrated our framework into the n -level hypergraph partitioner *KaHyPar*. A *flow-based refinement* is executed in $\log n$ levels of the multilevel hierarchy between each adjacent block in the quotient graph. The pairwise block scheduling refinement is implemented in rounds and terminates if no hypernode changed its block in a round. The sizes of the flow problems are chosen adaptively. If a *flow* computation on two blocks yields an improvement the flow problem size is increased, otherwise it is decreased. Additionally, we try to automatically balance the partition after *Max-Flow-Min-Cut* computation by iterating over each minimum (S, T) -cutset. In the remaining levels, where no flow is performed, the classical *FM* heuristic is used to improve the quality of a partition. An observation during implementation was that only a minority of the *Max-Flow-Min-Cut* computations leads to an improvement of quality. Therefore, we implement several speed-up heuristics which prevents the execution of additional pairwise *flow* refinements.

Our new quality configuration *KaHyPar-MF* produced on 95% of our benchmark instances better partitions than our old baseline configuration *KaHyPar-CA*. On average the solution quality is 2% better and only within a factor of 2 slower. In comparison with other state-of-the-art hypergraph partitioners, *KaHyPar-MF* produced on 70% of the benchmark instances the best-known partitions with a running time comparable to the direct k -way implementation of *hMetis*.

7.1. Future Work

Due to the novelty of the approach, there is a lot of potential in optimizing our basic framework. We made a trade-off between time and quality to obtain a *High-Quality Hypergraph Partitioner* which runs in reasonable time. The quality mainly depends on the number of flow executions through the multilevel hierarchy. The number of flow executions depends on the running time

of the flow algorithm and the size of the flow problem. Optimizing those two basic building blocks of the framework will allow us to achieve better quality in the same amount of time. The flow network of a hypergraph proposed by Lawler [30] has a bipartite structure. Because of this structural regularity, there might be other more specialized flow algorithms which run faster on these types of networks. Therefore, a useful work would be to evaluate many different maximum flow algorithms on our benchmark set. Further, one could investigate if it is possible to maintain the whole flow network over the multilevel hierarchy without explicitly setting up the flow network before each flow execution. Also, it would be interesting if information from previous flow calculations can be used to speed-up the current flow calculation. Pistorius [40] described an algorithm which implicitly executes EDMOND-KARP on a hypergraph using labels on the hypernodes. In our first version of the framework, we also used a similar technique and implicitly executes a flow algorithm on an implicit representation of the underlying network. During experiments, it turned out that the explicit representation was up to a factor of 2-3 faster than the implicit version. We encountered several reasons for that behavior:

- (i) Our flow network represents a subhypergraph of the original hypergraph. Iterating over the edges of a node means to iterate also over hypernodes which are not part of the flow problem and therefore have to be ignored.
- (ii) There are many different cases when we want to increase the flow along an *augmenting path*.
- (iii) Many labels have to be introduced which lead to a large number of main memory accesses.
- (iv) Also the implicit flow network is not flexible enough. Adding a new sparsifying technique would require with great certainty a reimplemention of the implicit flow network.

In Section 5.3 and 6.5 we show that with three simple speed up heuristics our *flow-based refinement framework* is up to a factor of 2 faster with comparable quality. Therefore, it would be beneficial to further increase the effectiveness ratio of the flow computation by introducing more heuristics.

It is also possible to further sparsify the flow network. Assume there exists two hypernodes v_1 and v_2 with $d(v_1) = 3$ and $d(v_2) = 4$. Further, $|I(v_1) \cap I(v_2)| = 3$ which means that in each hyperedge e where $v_1 \in e$ also $v_2 \in e$ and there exists one hyperedge e' where $v_2 \in e'$ and $v_1 \notin e'$. All hypernodes with $d(v) \leq 3$ are removed in our hybrid flow network. Consequently, we would remove v_1 and insert a clique between all incident hyperedges. However, v_2 is part of the flow network and induced $2d(v_2) = 8$ edges. Alternatively, we could remove v_2 and expand the clique between all hyperedges of $I(v_1)$ with e' . In that case, we have to insert an edge from each hyperedge in $I(v_1)$ to e' and vice versa. Since $|I(v_1)| = d(v_1) = 3$ only $2|I(v_1)| = 6$ edges are induced and we can remove one hypernode. In general, an expansion of a k -clique to a $(k + i)$ -clique induced ik edges from the k nodes already contained in the clique to the i new nodes and $i(k + 1 - 1)$ edges from the i new nodes to the k nodes in the clique. If we can remove a hypernode from the flow network by expanding a k -clique between hyperedge nodes to a $(k + i)$ -clique, it is beneficial if the following inequality holds

$$ik + i(k + i - 1) = i^2 + 2ki - i \leq 2(k + i)$$

The inequality is only satisfied for $i = 1$. In this case, we can exactly remove 2 edges and 1 node from the flow network. A possible algorithm could be to sort the hypernodes according to their degree and for each hypernode store a clique label which indicates between how many incident hyperedges already exist a clique. Afterwards, we iterate over the hypernodes and if we remove a hypernode, we have to update the clique label of all hypernodes in the intersection of the currently inserted clique. We iterate over the hypernodes until none of the hypernodes could be removed anymore. However, we didn't find an efficient implementation of the above-described algorithm. The algorithm requires a fast calculation between the intersection of

several hyperedges. An explicit construction of the intersection hypergraph would occupy too much memory.

References

- [1] Y. Akhremtsev, T. Heuer, P. Sanders, and S. Schlag. Engineering a direct k-way Hypergraph Partitioning Algorithm. In *2017 Proceedings of the 19th Workshop on Algorithm Engineering and Experiments (ALENEX)*, pages 28–42. SIAM, 2017.
- [2] C. J. Alpert. The ISPD98 Circuit Benchmark Suite. In *Proceedings of the 1998 International Symposium on Physical Design*, pages 80–85. ACM, 1998.
- [3] C. J. Alpert and A. B. Kahng. Recent Directions in Netlist Partitioning: a Survey. *Integration, the VLSI journal*, 19(1-2):1–81, 1995.
- [4] R. Andersen and K. J. Lang. An Algorithm for Improving Graph Partitions. In *Proceedings of the 19th annual ACM-SIAM symposium on Discrete algorithms*, pages 651–660. Society for Industrial and Applied Mathematics, 2008.
- [5] D. A. Bader, H. Meyerhenke, P. Sanders, and D. Wagner. *Graph Partitioning and Graph Clustering*, volume 588. American Mathematical Society, 2013.
- [6] A. Belov, M. Heule, D. Diepold, and M. Järvisalo. The Application and the Hard Combinatorial Benchmarks in Sat Competition 2014. *Proceedings of SAT Competition*, pages 81–82, 2014.
- [7] T. N. Bui and C. Jones. Finding Good Approximate Vertex and Edge Partitions is NP-Hard. *Information Processing Letters*, 42(3):153–159, 1992.
- [8] A. Buluç, H. Meyerhenke, I. Safro, P. Sanders, and C. Schulz. Recent Advances in Graph Partitioning. In *Algorithm Engineering*, pages 117–158. Springer, 2016.
- [9] U. V. Catalyurek and C. Aykanat. Hypergraph-Partitioning-based Decomposition for Parallel Sparse-Matrix Vector Multiplication. *IEEE Transactions on parallel and distributed systems*, 10(7):673–693, 1999.
- [10] B. V. Cherkassky. A Fast Algorithm for Computing Maximum Flow in a Network. *Collected Papers*, 3:90–96, 1994.
- [11] B. V. Cherkassky and A. V. Goldberg. On Implementing the Push-Relabel Method for the Maximum Flow Problem. *Algorithmica*, 19(4):390–410, 1997.
- [12] T. A. Davis and Y. Hu. The University of Florida Sparse Matrix Collection. *ACM Transactions on Mathematical Software (TOMS)*, 38(1):1, 2011.
- [13] U. Derigs and W. Meier. Implementing Goldberg’s Max-Flow-Algorithm — A Computational Investigation. *Mathematical Methods of Operations Research*, 33(6):383–403, 1989.
- [14] S. Dutt and W. Deng. Vlsi circuit partitioning by cluster-removal using iterative improvement techniques. In *Proceedings of the 1996 IEEE/ACM international conference on Computer-aided design*, pages 194–200. IEEE Computer Society, 1997.
- [15] J. Edmonds and R. M. Karp. Theoretical Improvements in Algorithmic Efficiency for Network Flow Problems. *Journal of the ACM (JACM)*, 19(2):248–264, 1972.
- [16] C. M. Fiduccia and R. M. Mattheyses. A Linear-Time Heuristic for improving Network Partitions. In *Papers on Twenty-five years of electronic design automation*, pages 241–247. ACM, 1988.
- [17] L. R. Ford and D. R. Fulkerson. Maximal Flow through a Network. *Canadian journal of Mathematics*, 8(3):399–404, 1956.
- [18] L. R. Ford Jr and D. R. Fulkerson. *Flows in Networks*. Princeton university press, 2015.
- [19] S. Fortunato. Community Detection in Graphs. *Physics reports*, 486(3):75–174, 2010.
- [20] A. V. Goldberg and R. E. Tarjan. A new Approach to the Maximum-Flow Problem. *Journal of the ACM (JACM)*, 35(4):921–940, 1988.

- [21] T. Heuer. *Engineering Initial Partitioning Algorithms for direct k-way Hypergraph Partitioning*. PhD thesis, Karlsruher Institut für Technologie (KIT), 2015.
- [22] T. Heuer and S. Schlag. Improving Coarsening Schemes for Hypergraph Partitioning by Exploiting Community Structure. In *LIPICS-Leibniz International Proceedings in Informatics*, volume 75. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2017.
- [23] M. Holtgrewe, P. Sanders, and C. Schulz. Engineering a scalable High Quality Graph Partitioner. In *Parallel & Distributed Processing (IPDPS), 2010 IEEE International Symposium on*, pages 1–12. IEEE, 2010.
- [24] T. C. Hu and K. Moerder. Multiterminal Flows in a Hypergraph. In T. Hu and E. Kuh, editors, *VLSI Circuit Layout: Theory and Design*, chapter 3, pages 87–93. IEEE Press, 1985.
- [25] A. B. Kahn. Topological Sorting of Large Networks. *Communications of the ACM*, 5(11):558–562, 1962.
- [26] G. Karypis, R. Aggarwal, V. Kumar, and S. Shekhar. Multilevel Hypergraph Partitioning: Applications in VLSI Domain. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 7(1):69–79, 1999.
- [27] G. Karypis and V. Kumar. Multilevel k-way Hypergraph Partitioning. *VLSI design*, 11(3):285–300, 2000.
- [28] B. Krishnamurthy. An Improved Min-Cut Algorithm for Partitioning VLSI Networks. *IEEE Transactions on Computers*, (5):438–446, 1984.
- [29] K. Lang and S. Rao. A Flow-Based Method for improving the Expansion or Conductance of Graph Cuts. In *IPCO*, volume 4, pages 325–337. Springer, 2004.
- [30] E. L. Lawler. Cutsets and Partitions of Hypergraphs. *Networks*, 3(3):275–285, 1973.
- [31] T. Lengauer. *Combinatorial Algorithms for Integrated Circuit Layout*. Springer Science & Business Media, 2012.
- [32] H. Liu and D. Wong. Network Flow Based Multi-way Partitioning with Area and Pin Constraints. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 17(1):50–59, 1998.
- [33] Z. Á. Mann and P. A. Papp. Formula Partitioning Revisited. 2014.
- [34] K. Menger. Zur Allgemeinen Kurventheorie. *Fundamenta Mathematicae*, 10(1):96–115, 1927.
- [35] M. E. Newman. Analysis of Weighted Networks. *Physical review E*, 70(5):056131, 2004.
- [36] V. Osipov and P. Sanders. n-Level Graph Partitioning. In *European Symposium on Algorithms*, pages 278–289. Springer, 2010.
- [37] D. A. Papa and I. L. Markov. Hypergraph Partitioning and Clustering, 2007.
- [38] S. B. Patkar, H. Sharma, and H. Narayanan. Efficient Network Flow Based Ratio-Cut Netlist Hypergraph Partitioning. *WSEAS Transactions on Circuits and Systems*, 3(1):47–53, 2004.
- [39] J.-C. Picard and M. Queyranne. On the Structure of all Minimum Cuts in a Network and Applications. *Combinatorial Optimization II*, pages 8–16, 1980.
- [40] J. Pistorius and M. Minoux. An Improved Direct Labeling Method for the Max–Flow Min–Cut Computation in Large Hypergraphs and Applications. *International Transactions in Operational Research*, 10(1):1–11, 2003.
- [41] L. A. Sanchis. Multiple-Way Network Partitioning. *IEEE Transactions on Computers*, 38(1):62–81, 1989.

-
- [42] P. Sanders and C. Schulz. Engineering Multilevel Graph Partitioning Algorithms. In *ESA*, volume 6942, pages 469–480. Springer, 2011.
 - [43] S. Schlag, V. Henne, T. Heuer, H. Meyerhenke, P. Sanders, and C. Schulz. k-way Hypergraph Partitioning via n-Level Recursive Bisection. In *2016 Proceedings of the Eighteenth Workshop on Algorithm Engineering and Experiments (ALENEX)*, pages 53–67. SIAM, 2016.
 - [44] D. D. Sleator and R. E. Tarjan. A Data Structure for Dynamic Trees. In *Proceedings of the thirteenth annual ACM symposium on Theory of computing*, pages 114–122. ACM, 1981.
 - [45] R. Tarjan. Depth-First Search and Linear Graph Algorithms. *SIAM journal on computing*, 1(2):146–160, 1972.
 - [46] N. Viswanathan, C. Alpert, C. Sze, Z. Li, and Y. Wei. The DAC 2012 Routability-Driven Placement Contest and Benchmark Suite. In *Proceedings of the 49th Annual Design Automation Conference*, pages 774–782. ACM, 2012.
 - [47] D. B. West et al. *Introduction to Graph Theory*, volume 2. Prentice hall Upper Saddle River, 2001.
 - [48] H. H. Yang and D. Wong. Efficient Network Flow Based Min-Cut Balanced Partitioning. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 15(12):1533–1540, 1996.
 - [49] Z. Zhao, L. Tao, and Y. Zhao. An Effective Algorithm for Multiway Hypergraph Partitioning. *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, 49(8):1079–1092, 2002.

A. Benchmark Instances

A.1. Parameter Tuning Benchmark Set

Type	Num	min V	Avg. V	max V	min E	Avg. E
ISPD98	5	32498	49049	69429	34826	52202
PRIMAL	5	53919	90467	163622	245440	414577
LITERAL	5	96430	141622	283720	140968	323388
DUAL	5	100384	297768	1070757	34317	85669
SPM	5	12328	34129	74104	12328	34129
Type	max E	Avg. e	Med. e	Avg. $d(v)$	Med. $d(v)$	Avg. $\frac{ E }{ V }$
ISPD98	75196	3.79	2	4.04	3.57	1.06
PRIMAL	629461	2.56	2.3	11.74	6.54	4.58
LITERAL	629461	2.56	2.3	5.85	3.25	2.28
DUAL	229544	8.05	6.03	2.32	2	0.29
SPM	74104	20.91	19.92	20.91	17.87	1

Table 5: Statistical summary of the parameter tuning instances.

A.2. Benchmark Subset

Type	Num	min V	Avg. V	max V	min E	Avg. E
DAC	5	522482	708389	917944	511685	697951
ISPD98	10	53395	110344	210613	60902	119535
PRIMAL	30	7729	141143	1613160	29194	632173
LITERAL	30	15458	281238	3226318	29194	632173
DUAL	30	29194	632173	6429816	7729	141143
SPM	60	11028	64765	1000005	4371	59589
Type	max E	Avg. e	Med. e	Avg. $d(v)$	Med. $d(v)$	Avg. $\frac{ E }{ V }$
DAC	898001	3.37	2	3.32	3.18	0.99
ISPD98	201920	3.87	2.08	4.2	3.67	1.08
PRIMAL	6429816	2.58	2.2	11.54	7.39	4.48
LITERAL	6429816	2.58	2.2	5.79	3.78	2.25
DUAL	1613160	11.54	7.39	2.58	2.2	0.22
SPM	1000005	16.25	12.95	14.95	12.58	0.92

Table 6: Statistical summary of the benchmark subset instances.

A.3. Full Benchmark Set

Type	Num	min V	Avg. V	max V	min E	Avg. E
DAC	10	522482	888090	1360217	511685	876629
ISPD98	18	12752	59801	210613	14111	64240
PRIMAL	92	7502	111371	1621762	28770	649991
LITERAL	92	15004	221981	3226318	28770	649991
DUAL	92	28770	649991	13378617	7502	111371
SPM	184	10000	56930	9845725	163	52709
Type	max E	Avg. e	Med. e	Avg. $d(v)$	Med. $d(v)$	Avg. $\frac{ E }{ V }$
DAC	1340418	3.41	2	3.37	3.27	0.99
ISPD98	201920	3.83	2.05	4.11	3.52	1.07
PRIMAL	13378617	2.74	2.31	16.01	8.12	5.84
LITERAL	13378617	2.74	2.31	8.03	3.65	2.93
DUAL	1621762	16.01	8.12	2.74	2.31	0.17
SPM	6920306	15.72	12.15	14.56	10.99	0.93

Table 7: Statistical summary of the full benchmark set instances.

A.4. Excluded Test Instances

Hypergraph	2	4	8	16	32	64	128
10pipe-q0-k.dual				△	△	△	○△
10pipe-q0-k.primal	□	□	□	□	□	□	□
11pipe-k.dual	△	○△	○△	○△	○△	○△	○△
11pipe-k				○	○	○	○
11pipe-k.primal	□	□	□	□	□	□	○□
11pipe-q0-k.dual					△	○△	○△
11pipe-q0-k.primal	□	□	□	□	□	□	□
9dlx-vliw-at-b-iq3.dual							△
9dlx-vliw-at-b-iq3.primal	□	□	□	□	□	□	□
9vliw-m-9stages-iq3-C1-bug7.dual	△	●○△	●○△	●○△	●○△	●○△	●○△
9vliw-m-9stages-iq3-C1-bug7	△	△	●○△	●○△	●○△	●○□△	●○□△
9vliw-m-9stages-iq3-C1-bug7.primal	△	△		△	○△	○△	○△
9vliw-m-9stages-iq3-C1-bug8.dual	△	●○△	●○△	●○△	●○△	●○△	●○△
9vliw-m-9stages-iq3-C1-bug8	△	△	●○△	●○△	●○△	●○□△	●○□△
9vliw-m-9stages-iq3-C1-bug8.primal	△	△		△	○△	○△	○△
blocks-blocks-37-1.130-NOTKNOWN.dual	○	●○	●○	●○	●○	●○	●○△
blocks-blocks-37-1.130-NOTKNOWN	□		□	□	□	□	□
blocks-blocks-37-1.130-NOTKNOWN.primal	□	□	□	□	□	□	□
E02F20.dual							○
E02F22.dual						○	○
openstacks-p30-3.085-SAT.primal	□	□	□	□	□	□	□
openstacks-sequencedstrips-nonadl-	□	□	□	□	□	□	□
nonnegated-os-sequencedstrips-p30-3.025-							
NOTKNOWN.primal							
openstacks-sequencedstrips-nonadl-	□	□	□	□	□	□	□
nonnegated-os-sequencedstrips-p30-3.085-							
SAT.primal							

A BENCHMARK INSTANCES

q-query-3-L100-coli.sat.dual							△
q-query-3-L150-coli.sat.dual							△
q-query-3-L200-coli.sat.dual							△
q-query-3-L80-coli.sat.dual							△
transport-transport-city-sequential-25nodes-							△
1000size-3degree-100mindistance-3trucks-							△
10packages-2008seed.030-NOTKNOWN.dual							△
transport-transport-city-sequential-	□						□
25nodes-1000size-3degree-100mindistance-							□
3trucks-10packages-2008seed.050-							□
NOTKNOWN.primal							□
velev-vliw-uns-2.0-uq5.dual			△	△	△	△	△
velev-vliw-uns-2.0-uq5.primal	□	□	□	□	□	□	□
velev-vliw-uns-4.0-9.dual				△	△	△	△
velev-vliw-uns-4.0-9.primal	□	□	□	□	□	□	□
192bit	□		□				
appu					○	○	
ESOC	□	□		□	○□	□	
human-gene2				○△	○△	○△	
IMDB			△	△	△	△	△
kron-g500-logn16	△	△	△	△	○△	○△	
Rucci1				□			
sls	□	□	□	○□	○□	○□	○□
Trec14							○

△ : KaHyPar-CA/KaHyPar-MF exceeded time limit
● : hMetis-R exceeded time limit
○ : hMetis-K exceeded time limit
□ : PaToH-Q memory allocation error

Table 9: Instances excluded from the full benchmark set evaluation.

B. Detailed Flow Network and Algorithm Evaluation

Instance	$ V' $	GOLDBERG-TARJAN				EDMOND-KARP			
		T_{Hybrid} $t[\text{ms}]$	T_G $t[\%]$	T_H $t[\%]$	T_L $t[\%]$	T_{Hybrid} $t[\%]$	T_G $t[\%]$	T_H $t[\%]$	T_L $t[\%]$
ALL	500	0.91	+2.24	+24.93	+29.35	-25.39	-24.3	-6.68	-11.53
	1000	1.95	+3.65	+26.19	+32.95	-13.99	-12.36	+10.81	+7.51
	5000	13.71	+8.63	+29.39	+43.11	+27.03	+35.33	+73.97	+86.31
	10000	30.54	+12.57	+36.15	+54.62	+47.93	+61.72	+100.41	+123.31
	25000	67.96	+23.36	+52.12	+87.8	+53.25	+77.85	+100.95	+138.8
DAC	500	0.34	-0.36	+30.14	+34.98	-37.61	-38.08	-23.12	-26.56
	1000	0.8	-1.7	+41.18	+47.43	-38.94	-41.19	-20.88	-22.17
	5000	5.2	+4.11	+46.02	+58.5	-21.35	-19.79	+12.55	+19.6
	10000	10.67	+3.2	+48.92	+66.83	-9.41	-6.44	+46.23	+63
	25000	31.43	+26.81	+186.2	+255.32	-23.53	-17.16	+25.16	+47.29
ISPD98	500	0.48	-0.58	+26.23	+28.54	-33.85	-34.5	-19.55	-20.14
	1000	1.11	-0.8	+32.35	+37.47	-29.32	-31.59	-11.91	-11.88
	5000	7.06	+6.65	+35.1	+49.35	-1.67	+1.64	+31.03	+41.91
	10000	16.33	+10.97	+42.54	+64.68	+18.38	+25.84	+75.19	+95.09
	25000	75.01	+26.26	+73.85	+132.06	+37.85	+56.79	+85.28	+124.01
DUAL	500	0.3	+12.37	+0.99	+13.6	-40.36	-34.35	-39.13	-37.67
	1000	0.6	+16.87	+0.83	+18.38	-40.93	-35.35	-39.47	-37.18
	5000	3.2	+37.54	+0.21	+37.78	-39.66	-23.77	-39.17	-24.01
	10000	5.78	+55.72	+1.21	+55.86	-34.01	-7.81	-33.3	-8
	25000	14.71	+105.19	+2.15	+105.88	-33.35	+17.43	-32.59	+17.28
PRIMAL	500	1.85	-0.73	+73.92	+76.03	+0.86	+0.17	+79.92	+63.57
	1000	3.9	+0.15	+77.48	+81.23	+33.02	+33.57	+160.43	+145.98
	5000	29.8	+0.84	+88.23	+96.71	+160	+162.28	+481.91	+510.71
	10000	45.94	+0.69	+109.75	+120.04	+195.68	+197.69	+487.6	+511.93
	25000	174.32	+0.21	+151.07	+159.04	+243.77	+248.81	+609.44	+648.46
LITERAL	500	0.86	+0.72	+63.65	+67.45	-16.1	-15.41	+35.63	+29.41
	1000	1.92	+1.64	+64.51	+71.46	+15.13	+17.07	+95.07	+90.72
	5000	12.31	+6.15	+76.65	+94.2	+59.04	+66.99	+216.7	+243.13
	10000	29.75	+8.55	+97.28	+115.37	+102.47	+117.45	+302.93	+363.17
	25000	64.4	+15.75	+128.34	+175.78	+126.59	+148.78	+286.31	+349.43
SPM	500	1.46	+0.35	+1.22	+2.47	-29.92	-30.42	-28.84	-34.57
	1000	3.09	+1.45	+1.14	+3.28	-23.32	-22.94	-22.17	-26.89
	5000	25.81	+1.79	+1.09	+3.26	+26.02	+28.55	+28.61	+27.43
	10000	74.81	+3.78	+2.48	+5.38	+45.86	+49.36	+48.77	+51.06
	25000	107.6	+6.67	+8.56	+12.07	+44.39	+48.88	+47.68	+52.96

Table 10: Running time comparison of maximum flow algorithms on different flow networks.

Note, all values in the table are in percentage relative to Goldberg-Tarjan on flow network T_{Hybrid} . In each line the fastest variant is marked bold.

C. Effectiveness Tests for Flow Configurations

To evaluate the effectiveness of our configurations presented in Section 6.4 we give each configuration the same amount of time to produce as many as possible partitions of a hypergraph H for a given k . We define $t_{H,k}$ which is the maximum partition time of a configuration to partition H in k blocks. If we execute a configuration on a hypergraph H for a given k and α' the time to produce as many as possible partitions is restricted by $3t_{H,k}$. We sum up the partition times during execution and if that sum plus the current average partition time would exceed $3t_{H,k}$ we perform the next run with a certain probability such that the expected running time is $3t_{H,k}$. The effectiveness tests were proposed by Sanders and Schulz [42]. The results of the tests mirrors our results of Section 6.4.

Config.	(+F,-M,-FM)	(+F,+M,-FM)	(+F,+M,+FM)
α'	Avg.[%]	Avg.[%]	Avg.[%]
1	-15.48	-15.22	0.14
2	-10.53	-10.11	0.36
4	-6	-5.13	0.66
8	-3.24	-1.68	1.24
16	-1.7	0.44	1.82
Ref.	(-F,-M,+FM)	6374.42	

Table 11: Table contains results of the effectiveness test for different configurations of our flow-based refinement framework for increasing α' . The quality in column *Avg.* is relative to our baseline configuration without the usage of flows.

D. Detailed Speedup Heuristic Evaluation

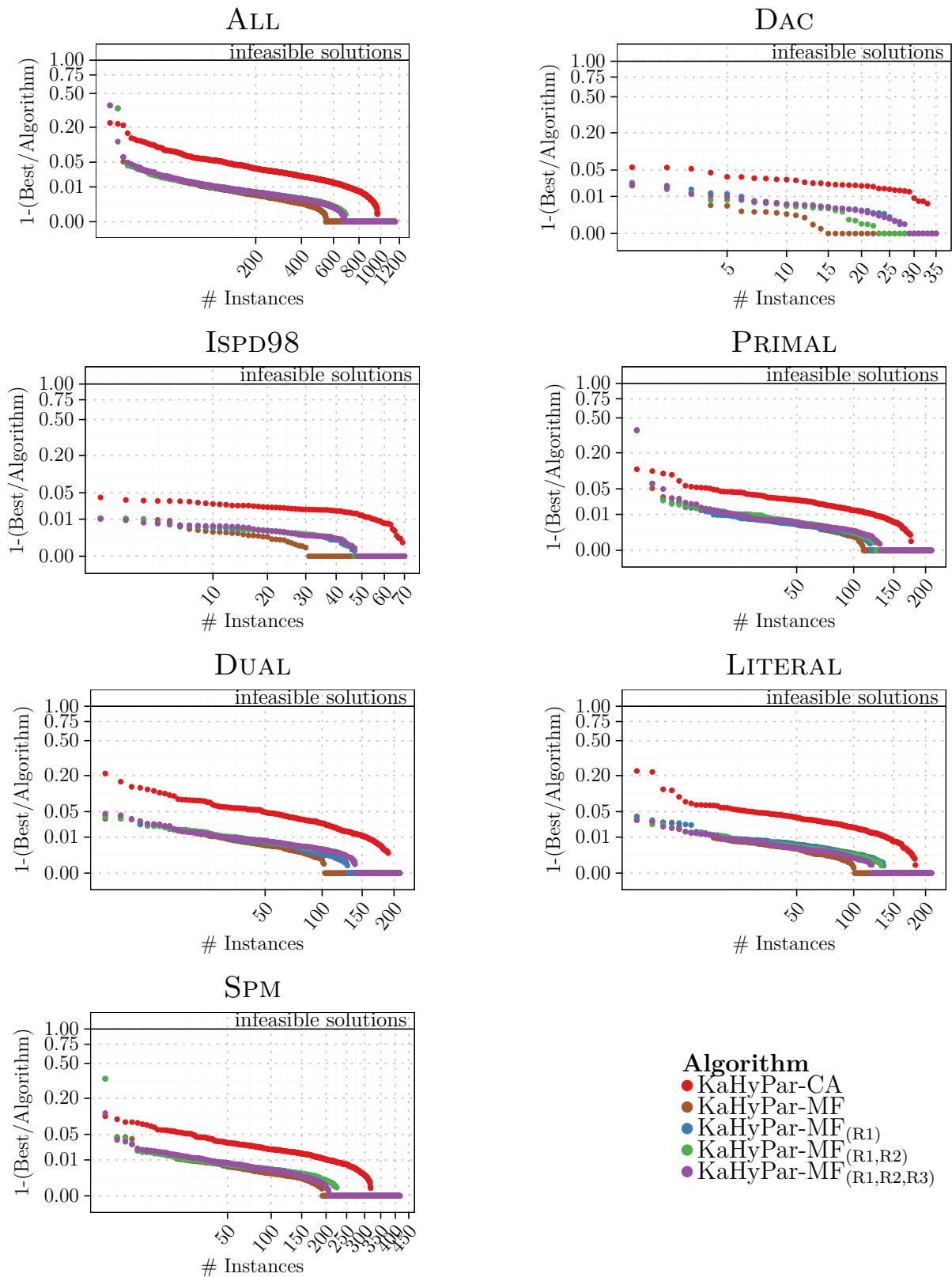


Figure 26: Min-Cut performance plots comparing KaHyPar-MF with KaHyPar-CA. The plots are explained in Section 6.2.

Partitioner	Running Time $t[s]$						
	ALL	DAC	ISPD98	PRIMAL	LITERAL	DUAL	SPM
KaHyPar-CA	29.26	343.4	21.57	36.44	56.49	58.75	11.31
KaHyPar-MF	81.54	699.18	75.97	114.67	185.56	143.93	28.74
KaHyPar-MF _(R1)	70.74	600.87	59.69	94.9	150.56	128.67	26.47
KaHyPar-MF _(R1,R2)	64.54	573.41	50.28	88.11	134.84	113.59	24.8
KaHyPar-MF _(R1,R2,R3)	56.88	526.86	43.32	74.76	116.79	101.76	22.31

Table 12: Comparing the average running time of KaHyPar-MF with KaHyPar-CA.

E. Detailed Comparison with other Systems

Partitioner	Average $\lambda - 1$						
	ALL	DAC	ISPD98	PRIMAL	LITERAL	DUAL	SPM
KaHyPar-MF	7819.11	17590.1	5671.37	15923.74	15844.61	3061.94	6165.74
KaHyPar-CA	2.03	2.47	1.72	1.69	2.25	2.71	1.75
hMetis-R	15.21	2.99	1.14	1.69	2.31	42.33	19.22
hMetis-K	14.71	7.78	0.9	3.66	8.77	27.66	19.09
PaToH-Q	8.98	12.86	7.41	11.72	12.81	7.96	6.37
PaToH-D	16.21	22.98	14.54	17.83	20.97	17.4	12.5

Table 13: Comparison of average ($\lambda - 1$) metric of KaHyPar-MF with KaHyPar-CA and other systems on different benchmark types. The results are in percentage relative to KaHyPar-MF.

Partitioner	Average $\lambda - 1$						
	$k = 2$	$k = 4$	$k = 8$	$k = 16$	$k = 32$	$k = 64$	$k = 128$
KaHyPar-MF	1064.06	3147.96	6062.8	9406	14756.03	21978.89	31820.94
KaHyPar-CA	1.73	2.06	2.36	2.28	2.11	1.9	1.73
hMetis-R	26.46	18.26	16.34	15.25	12.33	10.23	8.08
hMetis-K	26.86	17.19	15.18	15.06	11.29	9.83	8.1
PaToH-Q	11.1	8.5	8.57	9.49	8.87	8.6	7.7
PaToH-D	14.62	15.94	18.55	19.34	15.62	15.31	14.09

Table 14: Comparison of average ($\lambda - 1$) metric of KaHyPar-MF with KaHyPar-CA and other systems for different values of k . The results are in percentage relative to KaHyPar-MF.

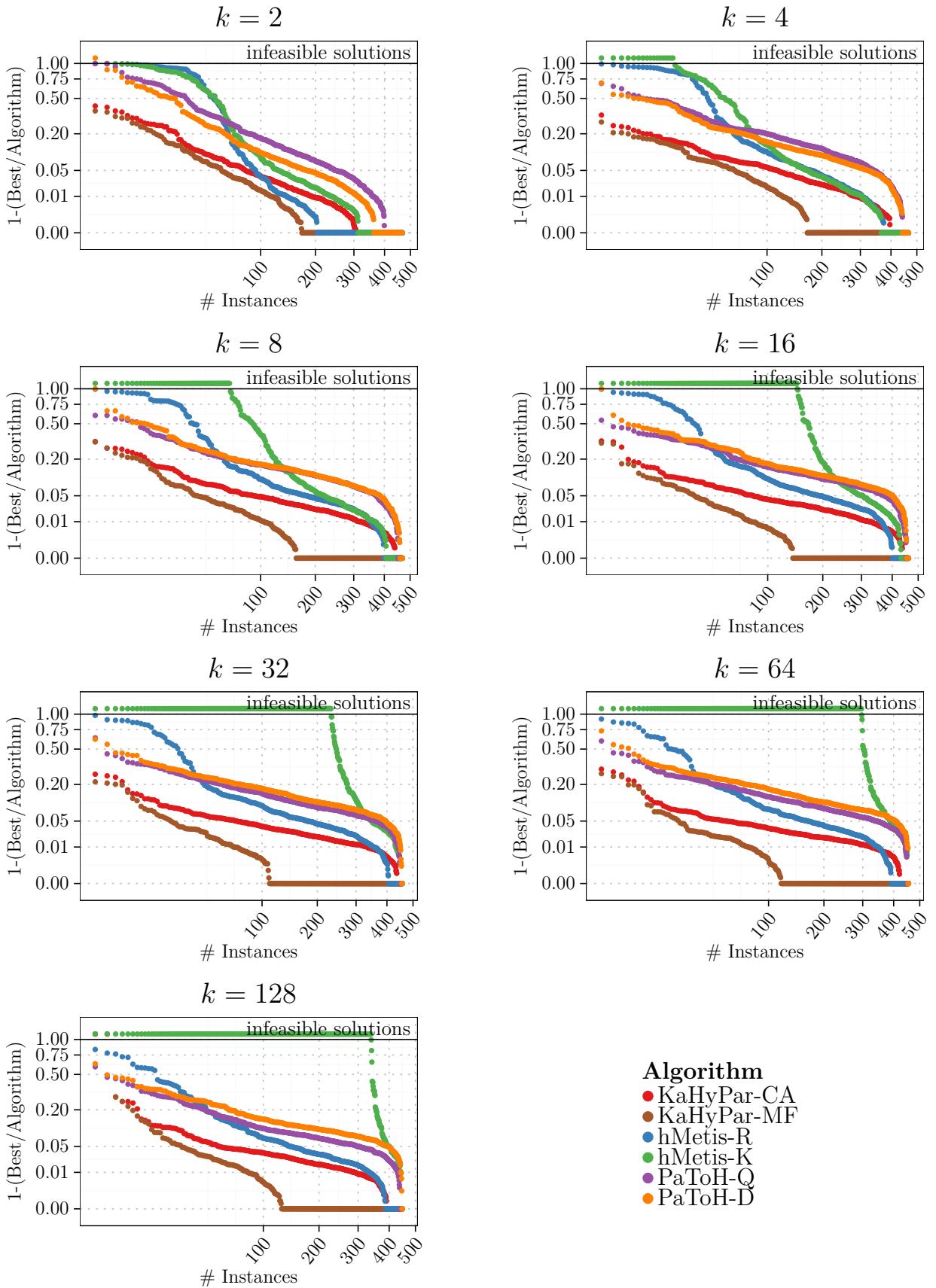


Figure 27: Min-Cut performance plots comparing KaHyPar-MF with KaHyPar-CA and other systems for different values of k .

Partitioner	Running Time $t[s]$						
	$k = 2$	$k = 4$	$k = 8$	$k = 16$	$k = 32$	$k = 64$	$k = 128$
KaHyPar-MF	22.13	38.51	55.04	67.83	85.75	108.97	128.04
KaHyPar-CA	12.68	17.16	23.88	31.01	41.69	57.35	76.61
hMetis-R	27.87	51.59	74.74	91.09	109.13	128.66	149.34
hMetis-K	25.47	32.27	42.5	53.41	74	109.12	152.92
PaToH-Q	1.93	3.61	5.44	7.01	8.4	10.06	11.44
PaToH-D	0.43	0.77	1.12	1.42	1.71	2.02	2.29

Table 15: Comparing the average running time of KaHyPar-MF with KaHyPar-CA and other systems for different values of k .