# High Quality Hypergraph Partitioning
# via Max-Flow-Min-Cut Computations
# Presentation-Script

Tobias Heuer

Karlsruhe Institute of Technology (KIT), Karlsruhe, Germany.

**Slide 1**  Hello and welcome, I am Tobias Heuer and today I am here to talk about my master thesis with the topic *High Quality Hypergraph Partitioning via Max-Flow-Min-Cut Computations*.

**Slide 2**  More precisely, my task was to develop a refinement algorithm to improve partitions of hypergraphs based on flow computations and integrate it into the $n$-level hypergraph partitioner *KaHyPar*. The well known *Max-Flow-Min-Cut* theorem establish an analogy between the minimum cut separating two vertices in a network and the maximum flow between them. Algorithms to solve the maximum flow problem are usually computationally expensive. A major goal of my work was to outperform the latest version of *KaHyPar* on most of our benchmark instances and simultaneously ensure that the running time remaining competitive.
As a motivation for my talk I will present right at the beginning our main contributions resp. our main experimental results.

- Our new approach outperform 5 different modern hypergraph partitioners on 73% of our large benchmark set.
- We imrpoved the quality of *KaHyPar* by 2.5%, while only incurring a slowdown by by a factor of 1.8.
- Our new version of *KaHyPar* has a comparable running time to *hMetis* and outperforms it on 84% of the instances.

Before we start to describe our framework, which leads to the presented results we need some fundamental definitions and related work.

**Slide 3**  Hypergraphs are generalizations of graphs where each hyperedge can connect more than two hypernodes. A hypergraph consists of a set of vertices, a set of hyperedges, where each hyperedge is a subset of the vertex set and a node and edge weight function. A hyperedge is also called net and vertex inside a net is called pin. The number of pins is the sum of the sizes of each net or the sum of the degree of each vertex.

**Slide 4**  The hypergraph partitioning problem is to partition the vertices of a hypergraph into $k$ non-empty disjoint blocks such that each block is smaller than $1 + \epsilon$ times the average block size. Simulaneously we want to minimize an objective function. In this talk we focus on the connectivity metric. Each hyperedge contributes $(\lambda - 1)$ times the edge weight to that function, where $\lambda$ is the number of blocks connected by net $e$. Consequently, the goal is to minimize the number of blocks inside a net $e$. **TODO 1:** *replace $\lambda$ with $\lambda(e)$ in slide*

**Slide 5**  Applications of hypergraph partitioning can be found in the area of VLSI Design. 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. A wire can connect more than two gates, therefore a hypergraph models

a circuit more accurately than a graph. **TODO 2:** *Ask Sebastian about the meaning of different applications*

**Slide 6** The hypergraph partitioning problem is NP-hard. The most common heuristic used in all modern hypergraph partitioners is the *multilevel paradigm*. First, a sequence of smaller hypergraphs is generated by contracting matchings or clusterings between vertices. If the hypergraph is small enough, expensive heuristics are used to calculate an *initial partition* of the hypergraph into $k$ blocks. Afterwards, the sequence of smaller hypergraphs is uncontracted in reverse order of contraction and the partition is projected to the next level finer hypergraph. After each uncontraction an *local search* algorithm is used to improve the hypergraph partition. Our *Max-Flow-Min-Cut* framework will work as a local search algorithm in this multilevel context.

**Slide 7** All modern hypergraph partitioners implements varations of the *FM* algorithm as local search algorithm. In general, this is a move-based heuristic that greedily move vertices between blocks based on local informations of incident nets. In this example we want to move the red hypernode from block 4 to block 3. The move reduces the cut of the partition by 1, because the light blue net is removed from cut. The reduction of the objective function when moving a node is also called gain. The *FM* heuristics moves a vertex with maximum gain in each step. **TODO 3:** *Move text* All modern hypergraph... *to the top*

**Slide 8** The *FM* heuristics is often critized, because it only incorparates local informations about the problem structure. If we want to move the red node in this example only the state of the colored nets is considered in the gain calculation. The quality of the final partition therefore heavily depends on the initial partition. In the multilevel context this problem is partially solved, because a move of a node corresponds to a move of a subset of the nodes in the original hypergraph, but the quality depends more on the quality of the coarsening phase.
Moreover, large hyperedges induce zero-gain moves. In this example, we have to perform many zero-gain moves until a single move of a vertex finally contributes to the gain function. In such situation the quality of the partition mainly depends on random choices made within the algorithm.
The motivation behind the usage of flow computations is that such an approach considers the global structure of the problem by finding a minimum-cut separating two vertices.

**Slide 9** Before we finally start with our main work, we need some basic terminology about network flows. Given a graph with a node and edge set and a capacity function $u$ and two vertices $s$ and $t$ which are also called source and sink. A flow is a function $f$ which satisfy the following two constraints:

– The flow on an edge has to be less or equal to the capacity of that edge
– And the amount of flow entering a node is the same as leaving the node.

The value of the flow is the amount of flow entering the sink resp. leaving the source. The maximum flow problem is to find a maximum flow $f$ from $s$ to $t$ such that no other flow function exists which value is greater than our maximum flow. An important concept to solve the maximum flow problem is the concept of the residual graph, which contains each edge with an residual capacity greater than zero. The residual capacity is the remaining amount of flow which we can send over an edge. On the reverse edge we say that the residual capacity is equal with the flow on the forward edge. An augmenting path is a path from $s$ to $t$ in the residual graph. One can show that $f$ is a maximum flow, if there exists no augmenting path.

**Slide 10**   We can calculate the corresponding minimum bipartition separating $s$ and $t$ by assinging each node reachable from the source $s$ in the residual graph to block 1 and all remaining to block 2.

**Slide 11**   Our flow-based reinfement framework is motivated by the results of Sanders and Schulz, who successfully integrate such an approach in their multilevel graph partitioner *KaFFPa*. We generalize many results of their work such that they applicable to hypergraph partitioning. In the following, we explain their basic ideas and highlight our contributions in different parts of the framework.

Our framework consists of four basic building blocks. We perform a flow-based refinement on two adjacent blocks of a $k$-way partition. Therefore, we need a block selection strategy, which is the first part of our framework. In the second part, we have to build a flow problem for our selected bipartition. Therefore, we have to add hypernodes to the flow problem and configure the source and sink set such that a *Max-Flow-Min-Cut* computations improves the cut of the corresponding bipartition. In the next step, we have to solve the flow problem. Therefore, we have to transform the subhypergraph into a corresponding flow network such that each mininimum-capacity cutset of the flow network corresponds to a minimum-weight cutset of the hypergraph. Finally, we use informations of the maximum flow function to automatically balance our bipartition according to our balanced constraints. One can show that one maximum flow computation has enough informations to enumerate all minimum-capacity cutsets.

We will now start to explain more in detail our block selection strategy.

**Slide 12**   Our block selection strategy is called *active block scheduling*. Each block of the partition has a state *active* or *inactive*. Initially all blocks are *active*. The algorithm works in rounds and starts by building the quotient graph of the $k$-way partition, which contains an edge between each adjacent pair of blocks. Then, we use our flow-based refinement algorithm to improve the quality of the bipartition induced by two adjacent *active* blocks. We visit the edges in random order. If a refinement yield an improvement, we mark the two blocks as *active* for the next round. In the first round, we found e.g., an improvement on block 1 and 2 and on all remaining pairs of adjacent blocks, there is no improvement. If the boundary of a block did not change, we mark them as inactive for the next round. In the next round, we only consider edges where one of the two blocks is *active*. In this example, no refinement yield an improvement in the second round. Consequently, all blocks become *inactive* and the algorithm terminates.

**Slide 13**   Once we select two adjacent blocks, we have to build a flow problem, which we use to improve the corresponding bipartition.

**Slide 14**  Our flow problem is build around the cut of the bipartition induced by the two selected blocks. We start a BFS, initialized with all vertices of a cut net of block 1. The BFS stops if the weight of the nodes would exceed a predefined upper bound. The upper bound is chosen in such a way that if all hypernodes of $B_1$ are assigned to block 2 the bipartition is still feasible according to our balanced contraint. Afterwards, we start a second BFS initialized with all vertices of a cut net of block 2. Let us consider all nets which are now partially contained in the flow problem. To ensure that non of the non-cut edges on the boundary becomes a cut edge we add the vertices of the nets, which are contained into flow problem, to the source resp. sink. Thus, we ensure that a *Max-Flow-Min-Cut* computation never increase the cut of the bipartition.

An extension of that approach is the *adaptive flow iteration* strategy. The size of the flow problem

depends now on $\varepsilon'$ rather than on $\varepsilon$. Then the algorithm works as follows: If a flow computation yields an improvement we double the flow problem size, where $\alpha'$ is a predifined upper bound. If a flow computation yields no improvement, we half the flow problem size. If $\alpha = 1$ yields no improvement the algorithm terminates.

**Slide 14** In the next step, we have to transform the subhypergraph extracted in the flow problem build step into a flow network.

**Slide 15** Therefore, we consider the bipartite graph representation of that hypergraph. Each hypernode is connected to its incident hyperedges via an undirected edge. Hu and Moerder introduce node capacities in this graph representation. Each hypernode as a weight equal to infinity and each hyperedge node has weight equal to its corresponding weight in the hypergraph. The observation is that a minimum-cutset separating two vertices of the hypergraph is equal with a minimum-weight vertex separator. We can use the known vertex separator transformation to obtain the corresponding flow network. Each node of the bipartite graph is splitted into an incoming and outgoing node, which are connected via an directed edge with the weight of that node. Lawler showed that it is not neccasary to split the hypernodes because their weight is infinity. The node containing all incoming edges is called incoming hyperedge node and the node containing all outgoing edges is called outgoing hyperedge node.

**Slide 16** If we go a step back and consider again the bipartite graph representation with its node capacities, we can make the following observation that a hypernode will never be part of a minimum-weight vertex separator, because their weight is infinity. Therefore, we remove all hypernodes by adding a clique between all incident hyperedges. Using the vertex separator transformation results in this network. A simple observation is that a hypernode induces 2 times the degree of that hypernode edges in the Lawler Network and the degree times the degree minus 1 edges in our network. We can reduce the number of nodes if the degree is smaller than 3. Therefore, we remove all low degree vertices in our flow network.

**Slide 17** Summing up the flow network construction, we can find a minimum-weight cutset in the hypergraph with the Lawler-Network. There also exists an other optimization due to Liu and Wong, which models each hyperedge of size 2 as a graph edge. In our network we remove each hypernode with a degree smaller than 3 by adding a clique between the corresponding hypernodes. Finally, we combine our and the Wong Network in a Hybrid Network to combine the advantages of both.

**Slide 18** We are now more familiar with the hypergraph flow network. Therefore, we consider again the flow problem modeling approach of Sanders and Schulz in *KaFFPa*. Each hypernode contained in a border hyperedges are added to the source resp. sink, which ensure that a non-cut edge did not become a cut edge after a *Max-Flow-Min-Cut* computation. This restricts the space of possible solutions, because non of the source resp. sinks vertices is moveable after a flow computation. We suggest an other flow problem modeling approach such that all vertices in the flow problem are moveable. Therefore, we extend the flow problem with all vertices contained in a border hyperedge and add those vertices to the source resp. sink. Now all vertices of our original flow problem are moveable, but the corresponding flow problem is significantly larger. Let us take a deeper look at the structure of our modeling approach. All incoming and outgoing edges of a source vertex have a capacity equal to infinity. Therefore, we can remove all previously added vertices and add all incident incoming hyperedge nodes to the source. We can do the same for the sink nodes and add all incident outgoing hyperedge nodes to the sink.