

Template \LaTeX Wiki von BAzubis für BAzubis

Projektarbeit 1 (T3_2000)

im Rahmen der Prüfung zum
Master of Science (B.Sc.)

des Studienganges Informatik

an der Dualen Hochschule Baden-Württemberg Karlsruhe

von

Vorname Nachname

Abgabedatum:	01. Februar 2025
Bearbeitungszeitraum:	01.10.2024 - 31.01.2025
Matrikelnummer, Kurs:	0000000, TINF15B1
Ausbildungsfirma:	SAP SE Dietmar-Hopp-Allee 16 69190 Walldorf, Deutschland
Betreuer der Ausbildungsfirma:	B-Vorname B-Nachname
Gutachter der Dualen Hochschule:	DH-Vorname DH-Nachname

Eidesstattliche Erklärung

Ich versichere hiermit, dass ich meine Projektarbeit 1 (T3_2000) mit dem Thema:

Template \LaTeX Wiki von BAzubis für BAzubis

gemäß § 5 der “Studien- und Prüfungsordnung DHBW Technik” vom 29. September 2017 selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Die Arbeit wurde bisher keiner anderen Prüfungsbehörde vorgelegt und auch nicht veröffentlicht.

Ich versichere zudem, dass die eingereichte elektronische Fassung mit der gedruckten Fassung übereinstimmt.

Karlsruhe, den August 5, 2025

Nachname, Vorname

Abstract

- English -

This is the starting point of the Abstract. For the final bachelor thesis, there must be an abstract included in your document. So, start now writing it in German and English. The abstract is a short summary with around 200 to 250 words.

Try to include in this abstract the main question of your work, the methods you used or the main results of your work.

Abstract

- Deutsch -

Dies ist der Beginn des Abstracts. Für die finale Bachelorarbeit musst du ein Abstract in deinem Dokument mit einbauen. So, schreibe es am besten jetzt in Deutsch und Englisch. Das Abstract ist eine kurze Zusammenfassung mit ca. 200 bis 250 Wörtern.

Versuche in das Abstract folgende Punkte aufzunehmen: Fragestellung der Arbeit, methodische Vorgehensweise oder die Hauptergebnisse deiner Arbeit.

Contents

Formelverzeichnis	VI
Abkürzungsverzeichnis	VII
List of Figures	VIII
List of Tables	IX
Quellcodeverzeichnis	X
1 Introduction	1
1.1 Motivation and Contribution	1
1.2 Structure	1
2 Preliminaries	2
2.1 Linear Optimization	2
2.2 Dantzig-Wolfe Decomposition	4
2.3 Graph Theory	5
2.4 Partition Refinement	6
2.5 Surprise and Entropy	9
2.6 Adjusted Rand Index	10
3 Related Work	12
4 Generic Column Generation (GCG)	13
4.1 Detection	14
4.2 Classifiers	16
4.3 Existing Detectors	21
4.4 Example	24
5 Tree Refinement	25
5.1 Information	25
5.2 Motivation	26
5.3 Classifiers	28
5.4 Variable Bounds	28
5.5 Relaxed MIPLIB types	29
5.6 Ordered Voting	30
5.7 Tree Refinement	31
5.8 Strategies	35

5.9	Scoring	39
6	Implementation	42
6.1	Architecture	43
6.2	Scoring	45
7	Evaluation	47
7.1	Setup	47
7.2	StrIPlib	48
7.3	Stuff	49
8	Notes	59
8.1	MIPLIB Constraint Types	59
8.2	Relaxed Constraint Types	60
8.3	Classifiers, Detectors	61
	Literaturverzeichnis	XI

Formelverzeichnis

A	mm ²	Fläche
D	mm	Werkstückdurchmesser
d_{\min}	mm	kleinster Schaftdurchmesser
L_1	mm	Länge des Werkstückes Nr. 1
	Grad	Freiwinkel
	Grad	Keilwinkel

Abkürzungsverzeichnis

GCG	Generic Column Generation
SCIP	SCIP
BaPCod	Branch-and-price Code
hMETIS	Hypergraph METIS
strIPlib	Structured Integer Program Library
MIP	Mixed Integer Programming

List of Figures

2.1	Solution space of a linear program highlighted in gray. Point highlighted in blue represent the solution space of the corresponding integer linear program.	3
2.2	A simple binpacking problem with 2 items and 1 bin represented as graph.	5
2.3	Entropy is measure of “surprise” and it increases with decreasing probability. On the left side, both colors are evenly distributed, so drawing either one is equally surprising. On the right side, drawing a red ball from the set of elements would be very surprising, because the probability is only $\frac{1}{12}$. But because this event is so unlikely, one does <i>not expect</i> to be surprised. As a result, the expected surprise - that is, the entropy - is low.	9
4.1	A simplified overview of the four major stages of solving a model with GCG.	13
4.2	A simplified overview of the detection process and its detection loop.	14
4.3	Visualization of the induced tree of propagated partial decompositions.	15
4.4	The graph of pair-wise Levenshtein weights for three capacity constraints. For $k = 1$, the edge between <i>capacity</i> ₃ and <i>capacity</i> ₁₂ vanishes, but because they is still a connecting path via. <i>capacity</i> ₁ , both constraints are assigned to the same class.	17
4.5	A constraint matrix with coefficients for each variable. Each constraint is assigned to a class corresponding to its number of non-zero entries.	18
4.6	Bin-Packing Model with items $\mathcal{I} = \{1, \dots, n\}$, item sizes $a_i \in \mathbb{Z}_{\geq 0}$, bins $\mathcal{J} = \{1, \dots, m\}$ and capacity C .	24
5.1	All parts of a model that contain useful information for semantic grouping of constraints and variables. Elements with a thick border are already used as a key concept in one of the existing detectors.	25
5.2	Test	31
7.1	The distribution of different categories of model within strIPlib. Most problems are part of “common” categories like Bin-Packing, Scheduling and Routing. The category “Others” includes e.g. Fantasy Football, Train Scheduling and different types for which only a small number of model files are available.	48

List of Tables

2.1	Contingency Table of partitions π and π' . Entry v_{ij} denotes the number of elements sets A_i and B_j have in common, i.e., $v_{ij} = A_i \cap B_j$	11
4.1	The structure of all 17 constraint types MIPLIB keeps track of.	20
4.2	For each classifier, the <i>cons class</i> detector will produce $2^k - 1$ new partial decompositions with k being the number of classes.	21
7.1	Consumer-grade components used to run all experiments.	47
8.1	A relaxed version of the constraint types MIPLIB uses.	60

Quellcodeverzeichnis

1 Introduction

1.1 Motivation and Contribution

1.2 Structure

2 Preliminaries

In this chapter, we provide sufficient background needed to understand the algorithmic details and ideas discussed in Chapter 5. First, we take a look at linear optimization as a general concept and discuss necessary basics such as Mixed Integer Programming (MIP). Afterwards, these concepts are used to take a more detailed look into the characteristics of the Dantzig-Wolfe decomposition, which represents a major component of the decomposition solver GCG. In addition, notations and basic definitions related to Graph Theory are introduced including the concept of *Partition Refinement*; an algorithmic approach used as a building block in Chapters 5 and 6.

2.1 Linear Optimization

Linear Optimization is a mathematical optimization technique used to determine the best possible values for a set of variables in a given model, whose constraints or requirements are represented by linear relationships. The goal is typically to maximize or minimize a objective function, subject to a set of equality and/or inequality constraints. Both objective and constraints must be linear.

If all variables are only allowed to take values from \mathbb{R}_{\geq}^n , i.e., only continuous values, then this optimization technique is referred to as *Linear Programming*. In a standard form, a linear programming problem with variable vector $\mathbf{x} \in \mathbb{R}^n$, constraint matrix $A \in \mathbb{R}^{m \times n}$, objective coefficients $c \in \mathbb{R}^n$ and right-hand side vector of the constraints $b \in \mathbb{R}^m$ can be expressed as follows:

$$z_{LP}^* = \min \quad c^T \mathbf{x} \tag{2.1}$$

$$\text{s.t.} \quad A \mathbf{x} \geq b \tag{2.2}$$

$$\mathbf{x} \geq 0 \tag{2.3}$$

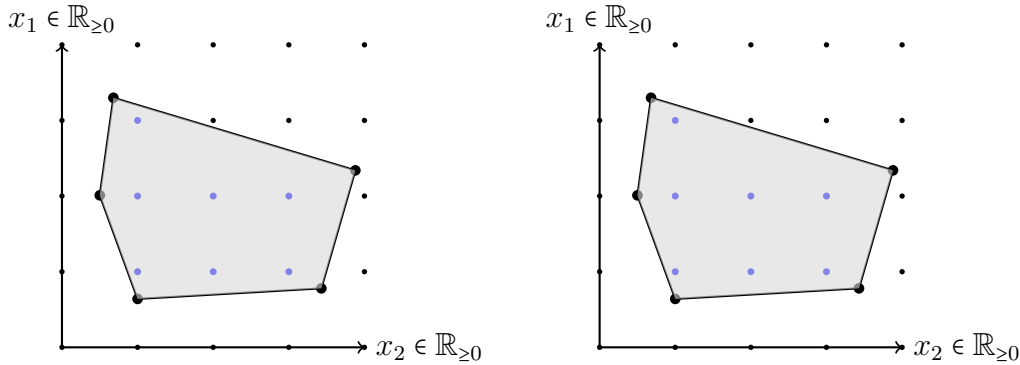


Figure 2.1: Solution space of a linear program highlighted in gray. Point highlighted in blue represent the solution space of the corresponding integer linear program.

Note that it is also possible to represent a set of equality $A'\mathbf{x} = b'$ by the two sets of inequalities $A'\mathbf{x} \geq b'$ and $A'\mathbf{x} \leq b'$. Without loss of generality, we assume optimization problems to always be minimization problems, unless explicitly stated otherwise. Constraints 2.2 specify a *convex* polytope over which the Objective Function 2.1 is optimized as shown in Figure 2.1. A solution vector $\mathbf{x} \in \mathbb{R}_{\geq 0}$ is called *feasible*, iff it satisfies both constraints 2.2 and 2.3. The linear program as a whole is called *feasible*, if there exists such a vector, otherwise it is considered *infeasible*. A feasible solution $\mathbf{x} \in \mathbb{R}_{\geq 0}$ is called *optimal* iff

$$c^T \mathbf{x} = \min\{c^T \mathbf{x} \mid \mathbf{x} \text{ is feasible}\} \iff \nexists \mathbf{y} \in \mathbb{R}_{\geq 0} : (\mathbf{y} \text{ feasible}) \wedge \underbrace{(c^T \mathbf{y} < c^T \mathbf{x})}_{\mathbf{y} \text{ is "better"}}$$

If there exists a feasible solution vector, but no optimal one, then the problem is called *unbounded*, as shown in Figure 2.1.

Linear programming is widely used in various fields such as operations research, economics, engineering, and logistics, due to its efficiency in solving large-scale real-world optimization problems. Algorithms such as the Simplex Method and Interior Point Methods are commonly used to solve LP problems efficiently. The simplex algorithm in particular is widely used in practice because of its efficiency on most problems, even though its worst-case complexity is exponential-time on certain families of problems depending on the chosen pivot-rule. However, on most problems the simplex algorithm only takes a polynomial number of steps to terminate. For more information about how the algorithms work and their mathematical details we refer to [1].

2.1.1 Mixed-Integer Programs

$$\begin{aligned}
 z_{LP}^* = \min \quad & c^T \mathbf{x} \quad + \quad d^T \mathbf{y} \\
 \text{s.t.} \quad & A \mathbf{x} \quad + \quad B \mathbf{y} \geq b \\
 & \mathbf{x} \quad \quad \quad \in \mathbb{R}_{\geq 0} \\
 & \quad \quad \quad \mathbf{y} \in \mathbb{Z}_{\geq 0}
 \end{aligned}$$

Even though linear programs are often sufficient to model a certain problem

2.2 Dantzig-Wolfe Decomposition

The *Dantzig-Wolfe Decomposition* is a

2.3 Graph Theory

Graphs are the fundamental data structure used in almost every aspect of computer science. This section will *not* introduce new concepts not already found in standard literature about graphs and related topics. We will mainly introduce the notation used for the following sections and chapters.

A *graph* is a tuple $G = (V, E)$ with $V \subseteq \{1, 2, \dots, n\}$ for some $n \in \mathbb{N}$ and $E \subseteq \{(u, v) \mid (u, v) \in V \times V\}$. The elements of set V are called *vertices* or *nodes*. The elements of set are ordered pairs called *directed edges*, *arcs* or simply *edges* which connect two vertices with each other. The set of outgoing neighbors of a specific vertex $v \in V$ is denoted $E(v) = \{(v, v') \mid v' \in V, vEv'\}$. The set of incoming neighbors $E^{-1}(v) = \{(v', v) \mid v' \in V, v'Ev\}$ is defined analogously. In this thesis, we will distinguish three types of graphs: directed, un-directed and bipartite, with directed being the assumed type if not stated explicitly.

For *un-directed* graphs, the edge relation E must be symmetric $\forall u, v \in V : uEv \rightarrow vEu$, i.e., if vertex u is connected to v or vice versa, then the corresponding back-edge must exist as well.

Bi-partite graphs are a special kind of graph class, where the set V can be represented with two sets $L, R \subseteq V$ such that $L \cap R = \emptyset$ and $E \subseteq \{(u, v) \mid u \in L, v \in R\}$. More informally, the vertex V can be split into two disjoint subsets $L, R \subseteq V$ such that no edges exists between vertices in each corresponding set. Bi-partite graphs are especially interesting, because the relationship between variables and the constraints they participate in can be encoded as such a graph as shown in Figure 2.2. This concept will be used in Chapter 4 to algorithmically detect various underlying structures.

Wording

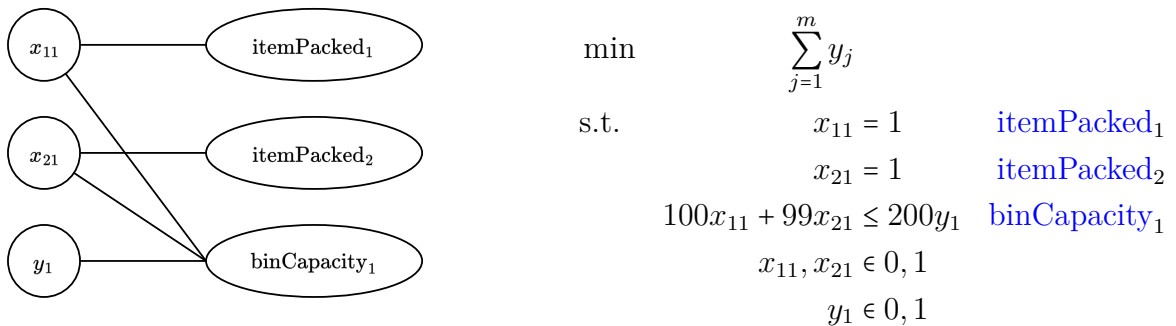


Figure 2.2: A simple binpacking problem with 2 items and 1 bin represented as graph.

2.4 Partition Refinement

Partition refinement is a fundamental concept in computer science, particularly relevant in fields such as automata theory [2], graph theory, and model checking [3]. A *partition* refers to a decomposition of a set U into disjoint, non-empty subsets $\{A_1, A_2, \dots, A_k\}$, called *cells* or *blocks*, such that:

$$\bigcup_{i=0}^k A_i = U \text{ and } \forall i \neq j : A_i \cap A_j = \emptyset$$

A partition $\pi = \{A_1, A_2, \dots, A_k\}$ of a set U is called a refinement of a partition $\pi' = \{B_1, B_2, \dots, B_m\}$, iff

$$\forall A_i \in \pi \exists B_j \in \pi' : A_i \subseteq B_j$$

As a special case, a partition is a refinement of itself. More informally, partition π' must reflect a “finer” classification of the elements than in π .

Partition refinement refers to an *iterative* process that refines a given initial partition of a set over the course of multiple iterations. Given a splitter-function $f : U \mapsto Q$ which maps every element of U to some element of an arbitrary target set Q , an initial partition π_{init} the goal is typically to find the coarsest partition $\pi_f = \{A_1, A_2, \dots, A_k\}$ of U such that the following two properties hold:

1. The partition π_f is a refinement of the initial partition π_{init}
2. $\forall A_i \in \pi_f \forall a, b \in P_i : f(a) = f(b)$, that is, the function f can intuitively be thought of a function expressing a certain “property” for each element. Elements with different properties cannot be part of the same cell and must be separated from each other during the refinement process.

Furthermore, the underlying problem structure to which partition refinement is applied, as well as the type of splitter function used, are not inherently restricted. In practice, however, many problems can be reformulated or encoded as graphs, where the function f captures a vertex property. For instance, in deterministic finite automaton (DFA) minimization, partition refinement is used to iteratively distinguish states by observing the equivalence classes of their transitions (Hopcroft’s algorithm): two states are grouped together only if, for every input symbol, their transitions lead into the same partition class; in graph isomorphism testing, it could encode vertex degrees or local neighborhood

structures; and in Markov decision processes (MDPs), f might reflect the expected reward or transition behavior. These encodings allow partition refinement to exploit structural symmetries and behavioral equivalences in a wide range of domains, especially if problems in that domain can be encoded as graphs. Further domains of application include Model Checking [3] and sorting algorithms [4].

For the purposes of this work, f will usually represent a function structurally similar to a *connection function* as it used in many graph automorphism packages. Given a graph $G = (V, E)$, then we define two types of connection function as follows:

$$f_{\text{count}}(v, X_{\text{ind}}) = |\{v' \in V \mid \forall (v, v') \in E, v' \in X_{\text{ind}}\}| \quad (2.4)$$

$$f_{\text{exists}}(v, X_{\text{ind}}) = \begin{cases} 1 & f_{\text{count}}(v, X_{\text{ind}}) \geq 1 \\ 0 & \text{else} \end{cases} \quad (2.5)$$

Functions 2.4 and 2.5

Algorithm 1 Partition refinement algorithm

while Test **do**
end while

Furthermore, if the underlying graph is bipartite and the splitter-function is expressing a vertex property, such as Function 2.4 or 2.5, then the partition refinement algorithm can be implemented more efficiently as highlighted in [5].

Algorithm 2 More efficient refinement, if graph $G = ((U, V), E)$ bipartite and $\forall v \in U : |E^{-1}(v)| \leq 1$ or $\forall v \in V : |E^{-1}(v)| \leq 1$. For the algorithm we assume the former.

Input: Initial partition $\Pi_{init} = \{A_1, A_2, \dots, A_k\}$, Bi-partite graph $G = ((L, R), E)$

Output: Coarsest stable partition

```

function REFINEFAST( $\Pi_{init}, G$ )
     $i \leftarrow -1$ 
     $\Pi_0 \leftarrow \Pi_{init}$ 
    for  $v \in R$  do
         $i \leftarrow i + 1$ 
         $\Pi_i \leftarrow \text{REFINE}(\Pi_{i-1}, E^{-1}(v))$ 
    end for
    return  $\Pi_i$ 
end function

```

2.5 Surprise and Entropy

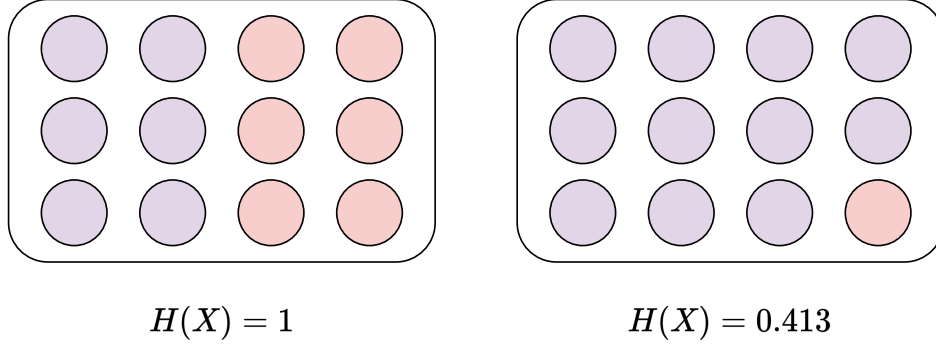


Figure 2.3: Entropy is measure of “surprise” and it increases with decreasing probability. On the left side, both colors are evenly distributed, so drawing either one is equally surprising. On the right side, drawing a red ball from the set of elements would be very surprising, because the probability is only $\frac{1}{12}$. But because this event is so unlikely, one does *not expect* to be surprised. As a result, the expected surprise - that is, the entropy - is low.

The *information value* or *surprisal* of an event E is defined as

$$I(E) = \log_b \left(\frac{1}{p(E)} \right) = -\log_b (p(E)) \quad (2.6)$$

It increases as the probability of the event $p(E)$ decreases. Intuitively, if the probability is close to 1, then one wouldn’t be surprised if this event actually occurred, so the surprisal is close to 0.

The *entropy*, or *expected surprise*, $H(X)$ of a discrete random variable X which takes values in the set \mathcal{X} is defined by equation 2.7 [6].

$$H(X) = \sum_{x \in \mathcal{X}} p(x) I(X) = - \sum_{x \in \mathcal{X}} p(x) \log_b p(x) \quad (2.7)$$

where $p(x) := \mathbb{P}[X = x]$.

If not specified any further, the base b of the logarithm is assumed to be 2. In chapter these concepts will be used to define a heuristic scoring system based on constraint names.

ref

2.6 Adjusted Rand Index

The *Rand Index* is a statistical measure used to compare two different partitions $\pi = \{A_1, A_2, \dots, A_k\}$, $\pi' = \{B_1, B_2, \dots, B_l\}$ of elements from the same set $U = \{1, 2, \dots, n\}$. Let $f_\pi(x) : U \mapsto \mathbb{N}$ be a function mapping an element $x \in U$ to the index of its cell in partition π . Function $f_{\pi'}$ is defined analogously. Furthermore, let

$$E_{\circ_1, \circ_2} = \{(x, y) \in U \times U \mid (f_\pi(x) \circ_1 f_\pi(y)) \wedge (f_{\pi'}(x) \circ_2 f_{\pi'}(y))\}$$

Intuitively, e.g. the set $E_{=, \neq}$ refers to the set of pairs $(x, y) \in U \times U$ of elements which are in the same cell in π , but in different cells in π' . Now we can define the *Rand Index* as follows:

$$\text{RI} = \frac{\overbrace{|E_{=, =}| + |E_{\neq, \neq}|}^{\text{Number of pairs for which } \pi, \pi' \text{ agree}}}{\underbrace{|E_{=, =}| + |E_{\neq, =}| + |E_{=, \neq}| + |E_{\neq, \neq}|}_{\text{Number of all pairs}}} = \frac{|E_{=, =}| + |E_{\neq, \neq}|}{\binom{n}{2}} \in [0, 1] \quad (2.8)$$

With appropriate data structures, e.g. a mapping between elements and cell index for each partition, Equation 2.8 can be evaluated in $O(n)$.

The *Adjusted Rand Index* is a chance-adjusted version of the regular *Rand Index* which accounts for similarities that might occur by random chance. It is one of the most popular measures for comparing partitions or clusters and can be computed by using Equation 2.9 [7].

$$\text{ARI} = \frac{\text{RI} - \text{Expected RI}}{\text{Max RI} - \text{Expected RI}} = \frac{\sum_{ij} v_{ij} \binom{v_{ij}}{2} - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] \div \binom{n}{2}}{\frac{1}{2} \left[\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right] - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] \div \binom{n}{2}} \quad (2.9)$$

The values v_{ij} , a_i and b_j are taken from the so called Contingency Table shown in Figure 2.1. We refer to [8] for a more detailed discussion on different similarity measures and their mathematical derivation.

π, π'	A_1	A_2	\dots	A_k	sums
B_1	v_{11}	v_{12}	\dots	v_{1k}	a_1
B_2	v_{21}	v_{22}	\dots	v_{2k}	a_2
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
B_l	v_{l1}	v_{l2}	\dots	v_{lk}	a_l
sums	b_1	b_2	\dots	b_k	

Table 2.1: Contingency Table of partitions π and π' . Entry v_{ij} denotes the number of elements sets A_i and B_j have in common, i.e., $v_{ij} = |A_i \cap B_j|$.

3 Related Work

4 Generic Column Generation (GCG)

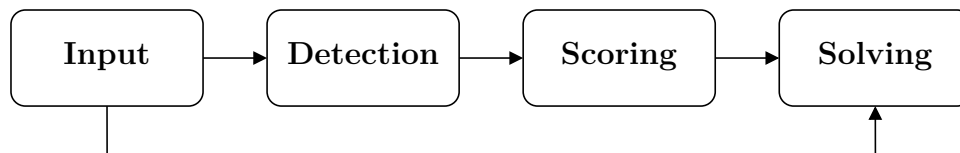


Figure 4.1: A simplified overview of the four major stages of solving a model with GCG.

In this chapter, we introduce Generic Column Generation (GCG), a decomposition solver which is based on the open-source MIP-Solver SCIP (SCIP) [9]. Readers already experienced with GCG and its capabilities may still find some details and observations interesting. For a given problem, GCG is able to perform an automatic Dantzig-Wolfe reformulation which is then solved using a branch-price-and-cut algorithm. Alternatively, GCG support a special *Benders-Mode* which reformulated the problem using Benders decomposition.

In contrast to other open-source solvers like BaPCod (Branch-and-price Code) [10] or commercial software such as *SAS* [11] which rely solely on user-provided decompositions, GCG is able to automatically detect different kinds of structures algorithmically, including but not limited to

- Single-Bordered structures
- Arrowhead structures using the third-party tool hMETIS (Hypergraph METIS) [12].
- Staircase structures

The solving process is divided into multiple consecutive stages as shown in Figure 4.1. Each stage will be explained in more detail in the following section as needed. The detection in particular aims to make GCG more accessible to a wider range of users which do not necessarily have the required theoretical background and practical experience to reformulate linear programs on their own. For more details about individual features and capabilities, we refer to the official documentation [13].

Entfernen
und auf
Kapitel
vorher
ref

Kurz
die 4
Schritte
aus Bild
erwäh-
nen und
einen
Satz

4.1 Detection

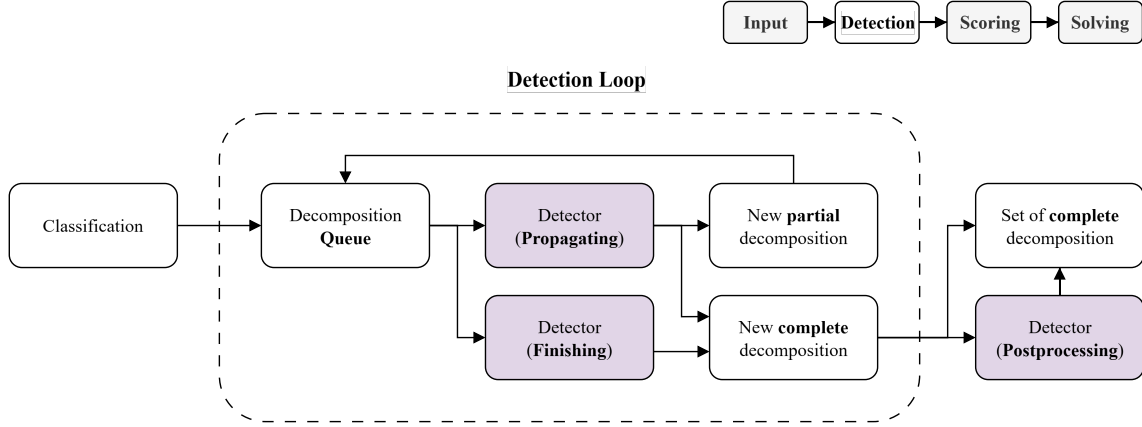


Figure 4.2: A simplified overview of the detection process and its detection loop.

As mentioned in the introduction to this chapter, one integral part and distinguishing feature of GCG is its detection framework. A simplified overview of the detection currently ¹ implemented in GCG is shown in Figure 4.2. For a more detailed visualization including additional information about how pre-solving is handled we refer to the official documentation [13]. The framework consists of two major parts:

1. A **classification** step, in which a set of classifiers is partitioning the constraints (and variables) according to a certain property, producing one partition each. The goal of this step is to detect different underlying structures of the constraint matrix, which can be used during the detection loop to make more informed decisions about which constraints to assign to which block or master. Important classifiers for the remainder of this thesis are discussed in more detail in Section 4.2.
2. The **detection loop**, which consists of a set of detectors which are responsible for assigning constraints either the master or to individual blocks. In round $n + 1$ a detector receives a *partial* decomposition, that is, a decomposition in which *not all* constraints are assigned yet, from round n as input and pushes a set of newly created (partial) decompositions to a queue. In case the user did not provide a partial decomposition as input in round 0, the loop is initialized with a decomposition in which no constraint is assigned yet.

¹GCG version 3.5, as of 2025-07-18.

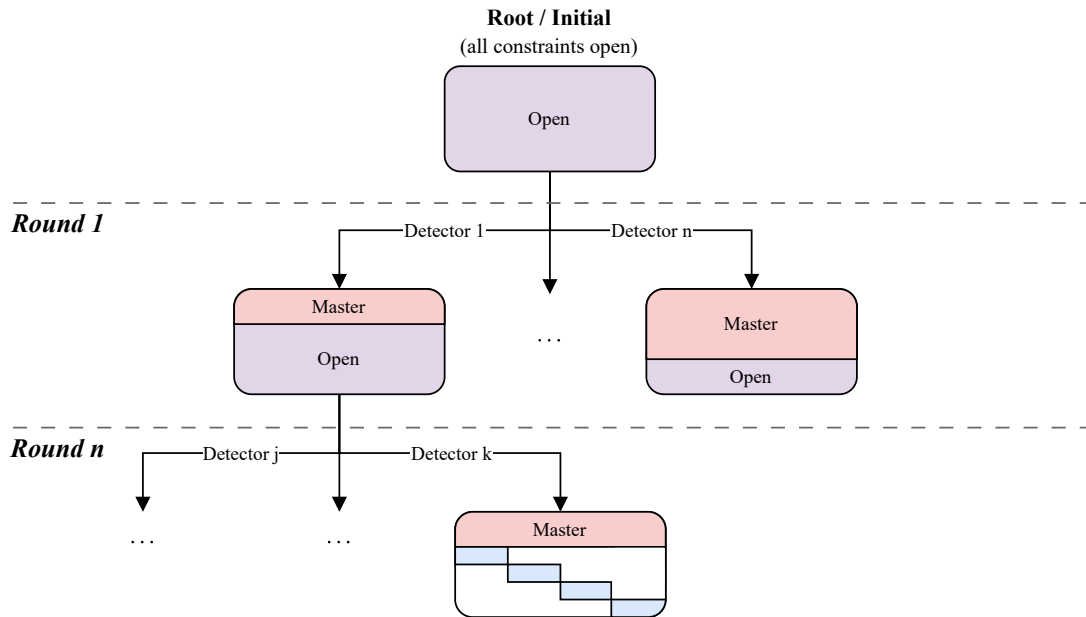


Figure 4.3: Visualization of the induced tree of propagated partial decompositions.

The concept of detecting structures in different rounds is visualized in Figure 4.3. Starting from a root decomposition in which all constraints are still unassigned or “open”, different detectors produce a set of new partial decomposition. Depending on the configuration, a detector is not allowed to work on a certain partial decomposition or its decedents twice. A very simple but concrete example of how such a tree might look like in practice can be found in Section 4.4.

Furthermore, if no detector found any new decomposition in round k , or k exceed the maximum number of rounds, the detection loop is stopped and all complete decomposition are collected, scored and exactly one is chosen for which the solving is started. The scoring and selection stage is of particular interest in practice, because the tree in Figure 4.3 might grow beyond thousand of decompositions, of which the best in terms of solving time or a different metric must be selected. Because the scoring of decompositions is not of major interest for *this* thesis, we refer to the official documentation for details [13].

Grammatik
Wort-
wahl

4.2 Classifiers

As mentioned in the introduction to this chapter, classifiers are responsible for detecting different underlying structures of the constraint matrix, which can be used during the detection loop to make more informed decisions about which constraints to assign to which block or master. Given a set of constraints $C = \{c_1, c_2, \dots, c_m\}$, classifiers can be seen as a *injective* function $f : C \mapsto \mathbb{Z}$, i.e., a function that assigns each constraint to exactly one number or *class*. Note that in GCG, classifiers are allowed to only classify a subset $C' \subseteq C$, leaving $C \setminus C'$ unassigned to any class ². In the current version of GCG, however, all classifiers always assign every constraint to some class.

4.2.1 Name Classifiers

The names of constraints and variables are, if provided, a strong indicator to which constraints or variables are related to each other. The names usually consist of two parts:

1. The semantic group name, such as “capacity” or “link” for e.g. a Bin-Packing model.
2. A *modifier*, which usually consists of numbers, capital letters or a combination of both. Typically, the modifier is separated from the semantic group name via. non alpha-numeric characters such as “_” or “#”.

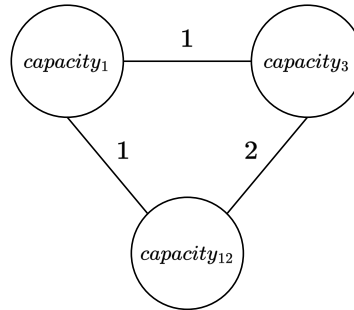
Constraints in the same group typically share similar names, with the *modifier* being the only differentiating factor. For example, in a Bin-Packing problem, capacity constraints such as “capacity_1”, “capacity_2”, ... usually vary only in the index indicating the bin. This similarity can be quantified using metrics like the *Levenshtein Distance*, which is the minimum number of single-character edits required to change one word into the other.

²When using GCG as a library, this can be checked via. `IndexPartition::isIndexClassified`.

Given a alphabet Σ , words $w, v \in \Sigma^*$ over that alphabet, then the *Levenshtein* distance between those two words can be computed as:

$$\text{lev}(w, v) = \begin{cases} |w| & \text{if } v = \epsilon \\ |v| & \text{if } w = \epsilon \\ \text{lev}(\text{prefix}(w), \text{prefix}(v)) & \text{if } \text{head}(w) = \text{head}(v) \\ 1 + \min \begin{cases} \text{lev}(\text{prefix}(w), v) \\ \text{lev}(w, \text{prefix}(v)) \\ \text{lev}(\text{prefix}(w), \text{prefix}(v)) \end{cases} & \text{otherwise} \end{cases} \quad (4.1)$$

Equation 4.1 can be computed in $O(|w| \cdot |v|)$ using a dynamic programming approach. Let $B = (\text{lev}(\text{name}(c_i), \text{name}(c_j)))_{1 \leq i, j \leq m}$ the pair-wise Levenshtein Distance between constraint names, $k \in \mathbb{N}$ the *connectivity* and $G = (V, E)$ with $V = \{c_1, c_2, \dots, c_m\}$, $E = \{\{u, v\} \mid u, v \in V, u \neq v, \text{lev}(\text{name}(u), \text{name}(v)) \leq k\}$. Furthermore, let $\text{reach}(v) = \{u \in V \mid (v, u) \in E\}$ be the set of reachable vertices from vertex $v \in V$. Then two constraints $c_i, c_j \in V$ are in the same class iff $c_j \in \text{reach}(c_i)$. A small example of this concept is shown in Figure 4.4.



5000
cons
limit

Figure 4.4: The graph of pair-wise Levenshtein weights for three capacity constraints. For $k = 1$, the edge between capacity_3 and capacity_{12} vanishes, but because they is still a connecting path via. capacity_1 , both constraints are assigned to the same class.

$$A = \begin{array}{c} \text{cons}_1 \\ \text{cons}_2 \\ \text{cons}_3 \\ \text{cons}_4 \end{array} \begin{array}{ccccc} x_1 & x_2 & x_3 & x_4 & \text{class} \\ \left(\begin{array}{cccc} 1 & 5 & -1 & -1 \\ 20 & 0 & 0 & 20 \\ 20 & 10 & 10 & 0 \\ 0 & 100 & -100 & 100 \end{array} \right) & \begin{array}{c} 4 \\ 2 \\ 3 \\ 3 \end{array} \end{array}$$

Figure 4.5: A constraint matrix with coefficients for each variable. Each constraint is assigned to a class corresponding to its number of non-zero entries.

4.2.2 Numeric Classifiers

Nonzero

The nonzero classifier classified constraints according to their number of non-zero variable coefficients as shown in Figure 4.5. Many types of models including Bin-Packing and Cutting-Stock consist only of constraint groups with a rather “stable” internal structure, i.e., the capacity constraint for each bin in model 4.4 consist of the same number of variables, because each constraint is just a sum over all items differing only in index for the respective bin. In general, constraint groups that are suited for this type of classifiers usually involve summations over fixed-sized sets (e.g. a set of items or bins) whose choice is not dependent on any quantified variable. Example for the latter include problems whose formulation is based on graphs and usually contains flow-conservation constraints shown in Equation 4.2.

Wrong
ref

$$\sum_{u \in E(v)} x_{uv} - \sum_{u \in E^{-1}(v)} x_{vu} = 0 \quad \forall v \in V \quad (4.2)$$

The amount of non-zeros in these constraint is entirely dependent on the number of outgoing and incoming edges for each vertex.

Objective

4.2.3 Type Classifiers

Classifiers based

SCIP Types

When using GCG as a library, the type of a variable or constraint can be retrieved via. `SCIPconsGetType(cons)` or `SCIPvarGetType(cons)` respectively. The former function is not provided by SCIP itself, but is implemented in GCG instead. The implementation compares the name of the handler the constraint is assigned to and compares it to a known list of constraint handlers. The list of supported handlers includes *Knapsack*, *Set Partitioning*, *Set Covering*, *Set Packing*, *Varbound* and *General*, in case no special structure was detected. Variables can be classified as *Integer*, *Binary* or *Continuous* ³.

the

Check
List

The clear downside of this classification is its important precondition. In order to use this feature properly and retrieve a meaningful type via. the two methods, pre-solving must have been executed prior to detection. When GCG reads the problem as e.g. an `.lp` file, all constraints are added as linear constraints to the underlying SCIP model. These constraints are usually “upgraded” if possible, that is, their structure is analyzed and assigned to the correct constraint handler during pre-solving. This is done to take advantage of properties only possessed by certain types of constraints, e.g. a solution to a set of Knapsack constraints *can* be computed more efficiently by using an algorithm based on dynamic programming. For more detailed information we refer to the official documentation [14].

Preliminary testing showed that it is not trivial to configure the pre-processing in such a way that *only* the upgrade mechanism is triggered and variables and constraints remain unchanged.

Add test
config to
appendix

³There are more types of variables in newer versions of SCIP such as *Implicit Integer*, but these three basic types are sufficient for the purpose of this discussion.

MIPLIB Constraint Types

Nr.	Type	Linear Constraint	Notes
1	Empty	\emptyset	-
2	Free	$-\infty \leq x \leq \infty$	No finite side.
3	Singleton	$a \leq x \leq b$	-
4	Aggregation	$ax + by = c$	-
5	Precedence	$ax - ay \leq b$	x, y have same type.
6	Variable Bound	$ax + by \leq c$	$x \in \{0, 1\}$
7	Set Partitioning	$\sum 1x_i = 1$	$\forall i : x_i \in \{0, 1\}$
8	Set Packing	$\sum 1x_i \leq 1$	$\forall i : x_i \in \{0, 1\}$
9	Set Covering	$\sum 1x_i \geq 1$	$\forall i : x_i \in \{0, 1\}$
10	Cardinality	$\sum 1x_i = b$	$\forall i : x_i \in \{0, 1\}, b \in \mathbb{N}_{\geq 2}$
11	Invariant Knapsack	$\sum 1x_i \leq b$	$\forall i : x_i \in \{0, 1\}, b \in \mathbb{N}_{\geq 2}$
12	Equation Knapsack	$\sum a_i x_i = 1$	$\forall i : x_i \in \{0, 1\}, b \in \mathbb{N}_{\geq 2}$
13	Bin Packing	$\sum a_i x_i + ay \leq a$	$\forall i : x_i, y \in \{0, 1\}, b \in \mathbb{N}_{\geq 2}$
14	Knapsack	$\sum a_i x_i \leq b$	$\forall i : x_i \in \{0, 1\}, b \in \mathbb{N}_{\geq 2}$
15	Integer Knapsack	$\sum a_i x_i \leq b$	$\forall i : x_i \in \mathbb{Z}, b \in \mathbb{N}$
16	Mixed Binary	$\sum a_i x_i + \sum p_j s_j \{ \leq, = \} b$	$\forall i : x_i \in \{0, 1\}, \forall j : s_j \text{ continuous}$
17	General Linear	$\sum a_i x_i \{ \leq, \geq, = \} b$	No special structure.

Table 4.1: The structure of all 17 constraint types MIPLIB keeps track of.

In order to get a more fine-grained classification based on constraint types. Because all types shown in Table 4.1 are deducible only from *local* information such as type of variables and right hand side coefficient, the types can be detected after one pass over the constraint matrix.

4.3 Existing Detectors

Nr.	Master	Open
1	(4.3)	(4.4)
2	(4.4)	(4.3)
3	(4.3), (4.4)	-

Table 4.2: For each classifier, the *cons class* detector will produce $2^k - 1$ new partial decompositions with k being the number of classes.

A

A

4.4 Example

$$\begin{array}{ll} \min & \sum_{j=1}^m y_j \\ \text{s.t.} & \sum_{j=1}^m x_{ij} = 1 \quad \forall i \in \mathcal{I} \end{array} \quad (4.3)$$

$$\begin{array}{ll} & \sum_{i=1}^n a_i x_{ij} \leq C y_j \quad \forall j \in \mathcal{J} \\ & x_{ij} \in 0, 1 \quad \forall i \in \mathcal{I}, \forall j \in \mathcal{J} \\ & y_j \in 0, 1 \quad \forall j \in \mathcal{J} \end{array} \quad (4.4)$$

Figure 4.6: Bin-Packing Model with items $\mathcal{I} = \{1, \dots, n\}$, item sizes $a_i \in \mathbb{Z}_{\geq 0}$, bins $\mathcal{J} = \{1, \dots, m\}$ and capacity C .



In order to illustrate the detection with a concrete example, we revisit the textbook Bin-Packing model shown in Figure 4.6. Constraints 4.3 enforce that every item is packed in exactly one bin, while inequalities 4.3 ensure that the capacity of each bin is respected if some item is packed in it. The objective is to minimize the number of bins.

Without pre-solving enabled, a classifier such as MIPLIB would assign constraints 4.3 and 4.4 to the classes *Set Partitioning* and *Bin-Packing* respectively. If unique, this classification is added to a list provided to the detection stage.

If no further classifications are found, the *cons class* detector will yield 3 new partial decompositions, first assigning constraints 4.3, then 4.4 and finally both 4.3 and 4.4 to the master. The constraint group not assigned to the master remains *open*.

During the next round of detection, a detector such as *Connected Base* will receive the partial decomposition with only the packing constraints assigned to the master as input. Here, the induced constraint adjacency graph of the $q \geq 0$ open capacity constraints consists of q isolated connected components, forming the desired block-diagonal structure.

5 Tree Refinement

With the existing capabilities of GCG presented in the previous chapter, we continue with the main contributions of this thesis:

- A new module which is integrated into the detection framework of GCG for reverse engineering semantic groupings of the original formulation.
- Additional auxiliary classifiers which implement constraint and variable classification based on information not currently used including examples of *when* they are crucial detecting semantics.

This chapter is divided into three main section:

1. A short summary about the available information we have access to.
2. What the motivation and goals are why and how we aim to process this information.
3. The concrete algorithm and its most integral parts.

Some concrete details about the implementation itself are not subject of the following sections, but are discussed in Chapter 6.

5.1 Information

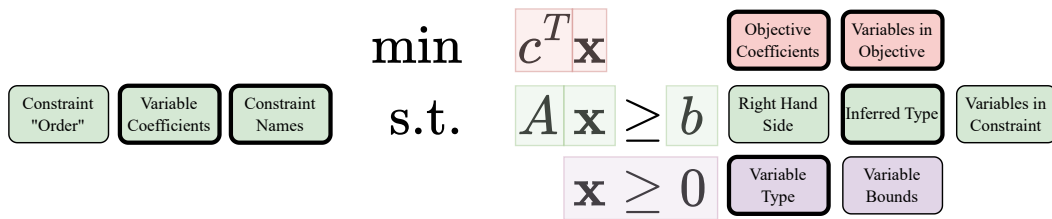


Figure 5.1: All parts of a model that contain useful information for semantic grouping of constraints and variables. Elements with a thick border are already used as a key concept in one of the existing detectors.

Before we present

5.2 Motivation

A

A

5.3 Classifiers

5.4 Variable Bounds

5.5 Relaxed MIPLIB types

5.6 Ordered Voting

5.7 Tree Refinement

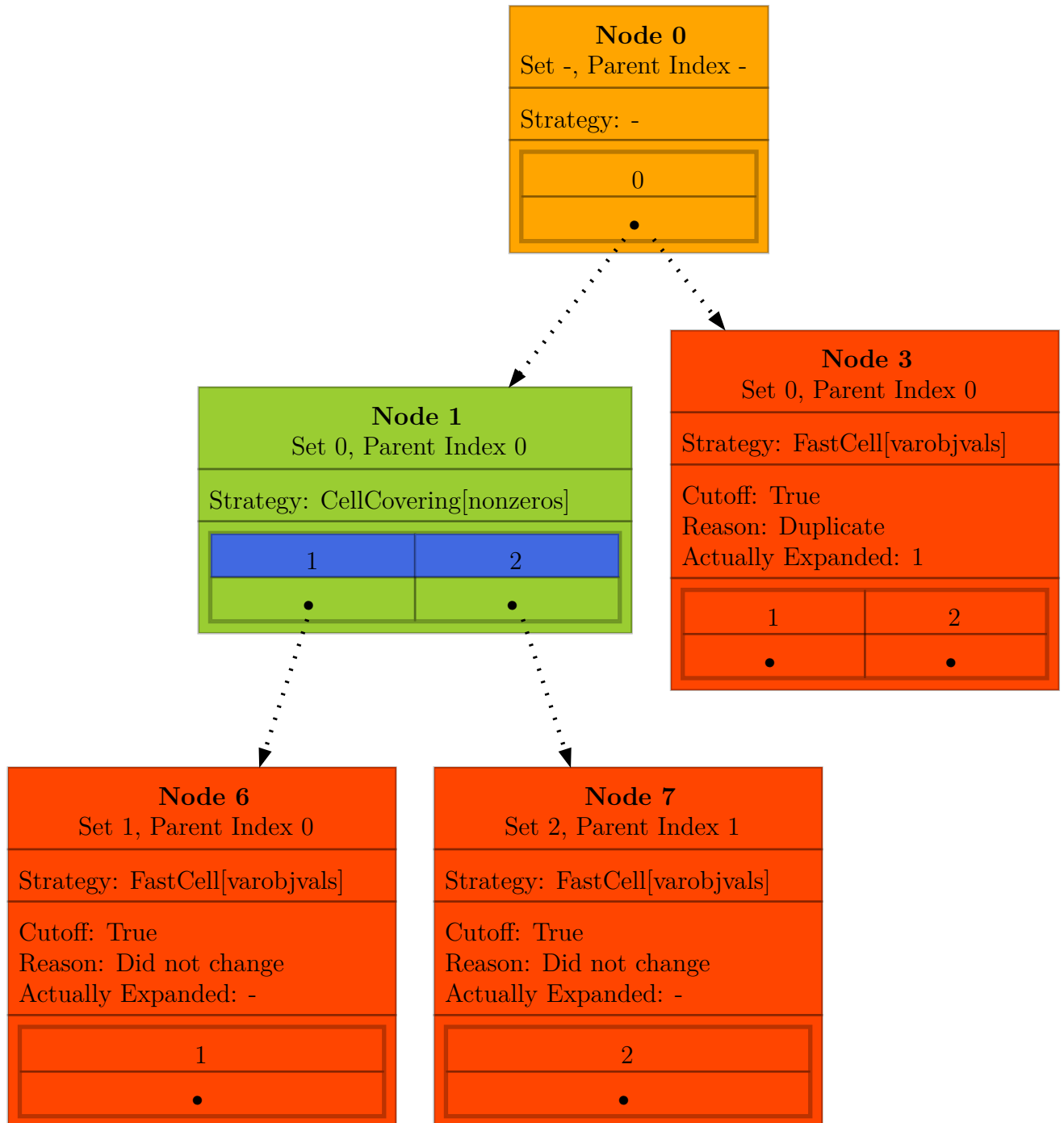


Figure 5.2: Test

A

A

A

5.8 Strategies

5.8.1 Slice

5.8.2 Fast

5.8.3 Recursive

5.8.4 Finishing

5.9 Scoring

5.9.1 Entropy Score

5.9.2 RAND Score

5.9.3 Connected Block Score

6 Implementation

6.1 Architecture

A

A

6.2 Scoring

6.2.1 Entropy Score

6.2.2 Connected Block Score

7 Evaluation

7.1 Setup

Type	Name	Metric
CPU	AMD Ryzen 3700X	3.8 GHz
RAM	-	16GB

Table 7.1: Consumer-grade components used to run all experiments.

All experiences were run on a system with components as specified in Table 7.1.

7.2 StrIPLib

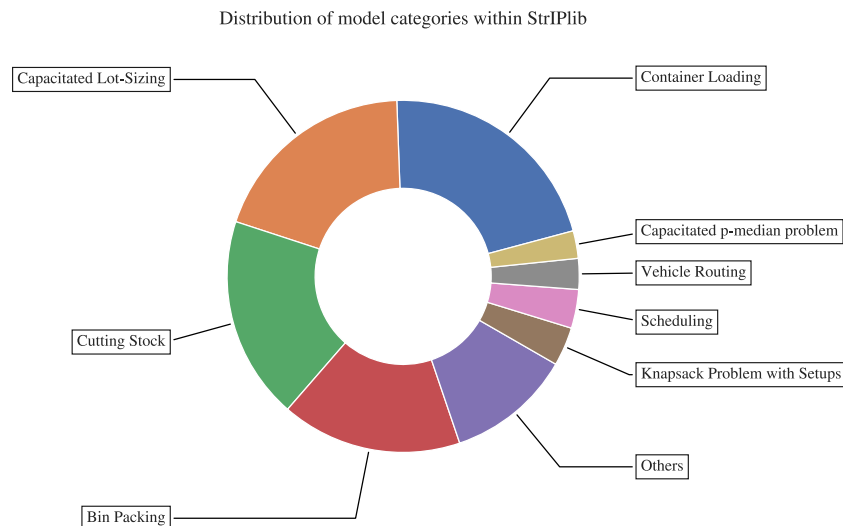


Figure 7.1: The distribution of different categories of model within strIPLib. Most problems are part of “common” categories like Bin-Packing, Scheduling and Routing. The category “Others” includes e.g. Fantasy Football, Train Scheduling and different types for which only a small number of model files are available.

The Structured Integer Program Library (strIPLib) is a collection of over 21000 mixed-integer programs with an exploitable structure such as Block-Diagonal and Staircase . All instances are assign to exactly one of the 33 main categories as highlighted in Figure 7.1. Each categories is further sub-divided into a number of smaller sub-categories, because each kind of problem can be modeled (e.g. three-index vs. four-index) *or* decomposed in a variety of different ways.

The number of available instances per category ranges from as low as 2 for Binary/Ternary Code Construction and up to ≈ 4700 for Container Loading, which makes the data-set in-balanced with respect to available models per main category. This is only of theoretical concern and is further discussed in section . Furthermore, the largest four categories account for $\approx 80\%$ of the total instance count with the remaining 20% distributed across 29 categories. One singular category “Haplotype Inference” with 40 instances is excluded from all tests, because the problem files are not readable by either GCG or SCIP. This behavior can be traced back to the used variable names in these models, which all contain the special character “^”.

ref

ref

Models
ohne
Namen
noch
ergänzen
(Prob-
lem-
beschrei

7.3 Stuff

A

A

A

A

A

A

A

A

A

8 Notes

8.1 MIPLIB Constraint Types

8.2 Relaxed Constraint Types

Nr.	Type	Linear Constraint	Notes
1	Empty	\emptyset	-
2	Free	$-\infty \leq x \leq \infty$	No finite side.
6	Variable Bound	$ax + by \leq c$	$x \in \{0, 1\}$
7	Set Partitioning	$\sum 1x_i = 1$	$\forall i : x_i \in \{0, 1\}$
8	Set Packing	$\sum 1x_i \leq 1$	$\forall i : x_i \in \{0, 1\}$
9	Set Covering	$\sum 1x_i \geq 1$	$\forall i : x_i \in \{0, 1\}$
10	Cardinality	$\sum 1x_i = b$	$\forall i : x_i \in \{0, 1\}, b \in \mathbb{N}_{\geq 2}$
10	Cardinality	$\sum 1x_i = b$	$\forall i : x_i \in \{0, 1\}, b \in \mathbb{N}_{\geq 2}$
11	Invariant Knapsack	$\sum 1x_i \leq b$	$\forall i : x_i \in \{0, 1\}, b \in \mathbb{N}_{\geq 2}$
12	Equation Knapsack	$\sum a_i x_i = 1$	$\forall i : x_i \in \{0, 1\}, b \in \mathbb{N}_{\geq 2}$
13	Bin Packing	$\sum a_i x_i + ay \leq a$	$\forall i : x_i, y \in \{0, 1\}, b \in \mathbb{N}_{\geq 2}$
14	Knapsack	$\sum a_i x_i \leq b$	$\forall i : x_i \in \{0, 1\}, b \in \mathbb{N}_{\geq 2}$
15	Integer Knapsack	$\sum a_i x_i \leq b$	$\forall i : x_i \in \mathbb{Z}, b \in \mathbb{N}$
16	Mixed Binary	$\sum a_i x_i + \sum p_j s_j \{ \leq, = \} b$	$\forall i : x_i \in \{0, 1\}, \forall j : s_j \text{ continuous}$
17	Mixed	$\sum a_i x_i \{ \leq, \geq \} b$	No special structure.
17	Mixed Equation	$\sum a_i x_i \{ = \} b$	No special structure.

Table 8.1: A relaxed version of the constraint types MIPLIB uses.

8.3 Classifiers, Detectors

Classifiers:

- VarTypes, ConsTypes (SCIP)
- MIPLIB Cons
- Name based / Levenstein
- Non-Zeroes
- Objective Values / Sign

Detectors:

- Cons Class, Var Class
- Connected Base
- HMETIS
- Dense Master Cons
- Detectors for single constraint types
- Staircase Heur
- Greedy
- Post Process

Literaturverzeichnis

- [1] *Branch & Price*.
- [2] “AN $n \log n$ ALGORITHM FOR MINIMIZING STATES IN A FINITE AUTOMATON”. In: Hopcroft, J. *Theory of Machines and Computations*. Elsevier, 1971, pp. 189–196. URL: <https://linkinghub.elsevier.com/retrieve/pii/B9780124177505500221> (visited on 07/20/2025).
- [3] Baier, C./ Katoen, J.-P. *Principles of Model Checking*. Cambridge, Mass: The MIT Press, 2008. 975 pp.
- [4] Mehlhorn, K./ Sanders, P. *Algorithms and Data Structures: The Basic Toolbox*. Berlin: Springer, 2008. 300 pp.
- [5] Salvagnin, D. “Detecting Semantic Groups in MIP Models”. In: *Integration of AI and OR Techniques in Constraint Programming*. Ed. by Quimper, C.-G. Vol. 9676. Cham: Springer International Publishing, 2016, pp. 329–341. URL: http://link.springer.com/10.1007/978-3-319-33954-2_24 (visited on 02/02/2025).
- [6] Cover, T. M./ Thomas, J. A. *Elements of Information Theory*. 2nd ed. Hoboken, N.J: Wiley-Interscience, 2006.
- [7] Sundqvist, M./ Chiquet, J./ Rigai, G. *Adjusting the Adjusted Rand Index – A Multinomial Story*. 11/17/2020. arXiv: 2011.08708 [stat]. URL: <http://arxiv.org/abs/2011.08708> (visited on 07/31/2025). Pre-published.
- [8] Warrens, M. J./ Van Der Hoef, H. “Understanding the Adjusted Rand Index and Other Partition Comparison Indices Based on Counting Object Pairs”. In: *Journal of Classification* 39.3 (11/2022), pp. 487–509. URL: <https://link.springer.com/10.1007/s00357-022-09413-z> (visited on 07/31/2025).
- [9] “Experiments with a Generic Dantzig-Wolfe Decomposition for Integer Programs”. In: Gamrath, G./ Lübbecke, M. E. *Lecture Notes in Computer Science*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 239–252. URL: http://link.springer.com/10.1007/978-3-642-13193-6_21 (visited on 07/08/2025).

- [10] Sadykov, R./ Vanderbeck, F. *BaPCod - a Generic Branch-and-Price Code*. Technical Report. Inria Bordeaux Sud-Ouest, 11/2021. URL: <https://inria.hal.science/hal-03340548>.
- [11] *SAS: Data and AI Solutions*. URL: https://www.sas.com/en_us/home.html (visited on 07/08/2025).
- [12] Karypis, G. et al. “Multilevel Hypergraph Partitioning: Application in VLSI Domain”. In: *Proceedings of the 34th Annual Conference on Design Automation Conference - DAC '97*. The 34th Annual Conference. Anaheim, California, United States: ACM Press, 1997, pp. 526–529. URL: <http://portal.acm.org/citation.cfm?doid=266021.266273> (visited on 07/08/2025).
- [13] *GCG*. URL: <https://gcg.or.rwth-aachen.de/> (visited on 07/08/2025).
- [14] *SCIP Doxygen Documentation: Overview*. URL: <https://www.scipopt.org/doc/html/> (visited on 07/08/2025).