

Thinning out facilities: a Benders decomposition approach for the uncapacitated facility location problem with separable convex costs

Matteo Fischetti¹, Ivana Ljubić², Markus Sinnl²

¹ Department of Information Engineering, University of Padua, Italy,
matteo.fischetti@unipd.it

² Department of Statistics and Operations Research, University of Vienna,
Austria, {ivana.ljubic,markus.sinnl}@univie.ac.at

March 27, 2015

Abstract

The Uncapacitated Facility Location (UFL) problem is one of the most famous and most studied problems in the Operations Research literature. Given a set of potential facility locations, and a set of customers, the goal is to find a subset of facility locations to open, and to allocate each customer to a single open facility, so that the facility opening plus customer allocation costs are minimized. For each customer, the allocation cost is assumed to be a linear or separable convex quadratic function. In this paper we “thin out” the classical models from the literature, and use generalized Benders cuts to replace a huge number of allocation variables by a small number of continuous variables that model the customer allocation cost directly. Although the idea of using Benders cuts for UFL is not new (at least, when linear costs are considered), it was apparently never computationally investigated in recent years. Instead, we show that the approach allows for a significant boost in the performance of a Mixed-Integer Programming solver, and report the optimal solution of a large set of previously unsolved benchmark. In particular, dramatic speedups are achieved for UFL’s with separable quadratic allocation costs—which turns out to be much easier than their linear counterpart when our approach is used.

The Uncapacitated Facility Location (UFL) problem is one of the most famous and most studied problems in the Operations Research literature. Given a set I of potential facility locations, and a set J of customers, the goal is to find a subset of facility locations to open, and to allocate each customer to a single open facility, so that the facility opening plus customer allocation costs are minimized. In its classical version, the allocation cost for each customer is assumed to be a linear function of the demand served by open facilities. The problem can be easily formulated as a compact Mixed-Integer Linear Program (MILP). In the last 50 years, two variants of this model (with aggregated and disaggregated constraints) have been traditionally used in the literature. Both variants however rely on a huge number of allocation variables, which is one of the main

reasons why UFL still imposes a challenge for modern general-purpose MILP solvers and Lagrangian relaxation techniques are generally preferred [4, 30, 37].

Recent experience with the Steiner Tree Problem [8] showed that some known MILP models can be “thinned out” by removing unnecessary variables and constraints, with a very significant performance boost. In this paper we apply the same approach to UFL, and use (generalized) Benders cuts to actually thin out the classical models. In the resulting formulation, the huge number of allocation variables is replaced by a linear number of continuous variables that model the customer allocation cost directly, or even by a single continuous variable. For UFL with linear costs, our Benders model is compact as it involves $|I| + |J|$ variables and $|J| \cdot (|I| + 1)$ constraints.

In addition to UFL, we also study the *quadratic UFL* (qUFL) problem in which customer demands are equal to one and allocation costs are proportional to the square of the fraction of the demand allocated to a given facility. This problem, also known as *separable convex quadratic UFL*, has been introduced in [18]. Due to its simple structure, qUFL has been a subject of intensive studies in the recent years. Many of the recently proposed methods for mixed integer nonlinear programming (MINLP) consider qUFL as an inevitable part of their computational studies, see e.g., [6, 7, 19, 20]. The importance of understanding computational techniques for qUFL is also confirmed by its presence in the Conic Benchmark library (CBLIB [13]), which is a library of the most relevant and challenging benchmark problems for conic optimization.

Although the idea of using Benders cuts for UFL with linear costs is definitely not new and can be considered folklore, to the best of our knowledge it was not computationally investigated in recent years. The outcome of our research is twofold:

1. On the one hand, the “thinning out” approach allows for a very significant boost in the performance of the MILP solver, as we were able to solve to proven optimality 7 previously unsolved benchmark instances for UFL, and to improve the best-known heuristic value for 22 additional instances. These instances were out of reach for any state-of-the-art MILP solver, as the underlying models would involve tens of millions of variables and constraints.
2. On the other hand, speedups of 4 orders of magnitude or more are reported for qUFL with respect to previous solution methods—to our surprise, qUFL turned out to be much easier than its linear counterpart when our approach is properly implemented.

The paper is organized as follows. The classical UFL models are reported in Section 1 for both the linear and separable quadratic cost cases. Nonlinear Benders cuts for convex optimization are reviewed in Section 2. Section 3 outlines our overall solution scheme, with a discussion of actual implementation issues that played a fundamental role for the effectiveness of our final solution algorithm. Computational results for UFL with both linear and separable quadratic

allocation costs are given in Section 4, while concluding remarks and possible future works are finally addressed in Section 5.

1 MIP models

Let I be the index set of facility locations ($|I| = n$), let $f_i \geq 0$ be opening costs for each facility $i \in I$, and let J be the index set of customers ($|J| = m$) with allocation costs $c_{ij} \geq 0$ defined for each pair $(i, j) \in I \times J$. We will assume without loss of generality that each customer can be allocated to every facility (if this is not the case, we will assume $c_{ij} = \infty$). In the traditional compact MILP formulation for UFL, $n + m \cdot n$ variables are used to model the problem. For each $i \in I$, binary variable y_i is set to one if facility i is open, and to zero, otherwise. For each $i \in I$ and $j \in J$, allocation variable x_{ij} is set to one, if customer j is served by facility i , and to zero otherwise.

Linear case

The classical UFL model for the linear case then reads

$$\min \sum_{i \in I} f_i y_i + \sum_{i \in I} \sum_{j \in J} c_{ij} x_{ij} \quad (1)$$

$$\text{s.t. } \sum_{i \in I} x_{ij} = 1 \quad \forall j \in J \quad (2)$$

$$x_{ij} \leq y_i \quad \forall i \in I, j \in J \quad (3)$$

$$x_{ij} \geq 0 \quad \forall i \in I, j \in J$$

$$y_i \in \{0, 1\} \quad \forall i \in I$$

The objective is to minimize the sum of facility opening costs, plus the customer allocation costs. Constraints (2) make sure that every customer is assigned to exactly one facility, and capacity constraints (3) make sure that allocation to a facility i is only possible if this facility is open. Note that the integrality condition on variables x_{ij} is redundant, i.e., for any integer y^* each customer j will set $x_{ij}^* = 1$ for the closest facility i with $y_i^* = 1$.

The model shown above is known as the disaggregated formulation, whereas in its aggregated counterpart $m \cdot n$ constraints of type (3) are replaced by n weaker constraints $\sum_{j \in J} x_{ij} \leq m \cdot y_i$, for all $i \in I$.

There is a large body of work available in the existing literature on exact and heuristic approaches for UFL. Recent survey articles focus on applications of facility location in design of distribution systems [24], or supply chain management [35], and a more general survey on UFL is given in [38]. Approaches applied to UFL range from approximation algorithms [32], over (semi-) Lagrangian relaxations [4], and metaheuristics (see, e.g., a survey in [3]), to branch-and-bound based algorithms in which lower bounds are calculated using dual ascent procedures (see, e.g. the most recent one in [31]). The state-of-the-art exact algorithm for UFL is given in [37]: the algorithm is based on a message passing approach

in which a metaheuristic calculates the upper bounds and passes the information to a branch-and-bound algorithm in which lower bounds are obtained by a Lagrangian relaxation solved using a bundle method. A more comprehensive overview of the most relevant recent literature can also be found in [37].

Separable convex quadratic case

UFL with separable convex quadratic allocation costs (denoted as qUFL in the following), has been introduced in [18]. In this version of the problem, one assumes $c_{ij} > 0$ for all $i \in I$ and $j \in J$, and the allocation costs are proportional to the square of a customer's demand served by an open facility. More precisely, the objective function (1) is replaced by:

$$\min \sum_{i \in I} f_i y_i + \sum_{i \in I} \sum_{j \in J} c_{ij} x_{ij}^2 \quad (4)$$

Contrarily to the linear case, each variable x_{ij} will now represent the fraction of demand of customer j served by facility i (in general, these values will no longer assume integer values in the optimal solutions).

Given a binary vector y^* , there is a simple closed formula to compute optimal allocation values for each customer. More specifically, for each customer $j \in J$ an optimal solution will set $x_{ij}^* = 0$ for all i with $y_i^* = 0$, and $x_{ij}^* = \delta_j^*/c_{ij}$ for all i with $y_i^* = 1$, where the normalization factor $\delta_j^* = 1/\sum_{i \in I: y_i^*=1} (1/c_{ij})$ guarantees the fulfillment of (2); see e.g. [18] for details.

Sophisticated linearization and bounding techniques for qUFL have been proposed in [18]. More recently, a very tight convex (second-order cone) MIP formulation has been given in [19], namely:

$$\min \sum_{i \in I} f_i y_i + \sum_{i \in I} \sum_{j \in J} c_{ij} z_{ij} \quad (5)$$

$$\text{s.t. } \sum_{i \in I} x_{ij} = 1 \quad \forall j \in J \quad (6)$$

$$x_{ij} \leq y_i \quad \forall i \in I, j \in J \quad (7)$$

$$x_{ij}^2 \leq z_{ij} y_i \quad \forall i \in I, j \in J \quad (8)$$

$$x_{ij} \geq 0 \quad \forall i \in I, j \in J \quad (9)$$

$$z_{ij} \geq 0 \quad \forall i \in I, j \in J \quad (10)$$

$$y_i \in \{0, 1\} \quad \forall i \in I \quad (11)$$

where (8) are second-order cone (hence convex) constraints as the right-hand-side is the product of nonnegative variables. The model above is called *perspective reformulation* (see, e.g., [12]) as it strengthens the obvious definition of the z_{ij} variables through constraints $x_{ij}^2 \leq z_{ij} y_i$, by replacing the left-hand side convex function x_{ij}^2 with its perspective defined by $y_i(x_{ij}/y_i)^2$ if $y_i > 0$, and zero if $y_i = 0$; see again [19] for details.

Computational experience reported in [19] shows that continuous relaxation of model (5)-(11), though very large, produces much tighter lower bounds than its counterpart based on (4), so we used it in our study.

2 Benders decomposition for convex optimization

Despite this broad set of solution approaches for UFL, it seems that Benders-like decomposition methods have been neglected so far—at least from a computational viewpoint. As already mentioned, the aim of the present paper is to close this gap and assess the computational performance of Benders decomposition for UFL and its quadratic counterpart with separable convex objective function.

Our overall framework works as follows: as x_{ij} variables are a bottleneck for MIP solvers, we just remove them from the model, and introduce in the objective function a new set of continuous variables w_j representing the allocation-cost for all $j \in J$. The resulting *master problem* is then given by

$$\min \sum_{i \in I} f_i y_i + \sum_{j \in J} w_j \quad (12)$$

$$\text{s.t. } \sum_{i \in I} y_i \geq 2 \quad (13)$$

$$w_j \geq \Phi_j(y) \quad \forall j \in J \quad (14)$$

$$y_i \in \{0, 1\} \quad \forall i \in I \quad (15)$$

where the convex function $\Phi_j(y)$ appearing in (14) gives the minimum allocation cost for customer j for any given (possibly noninteger) point $y \in [0, 1]^I$ with $\sum_{i \in I} y_i \geq 2$ (which is the only case of interest for branch-and-cut separation).

Note that we require to open, at least, 2 facilities in (13), as the single-facility case can be easily handled in a preprocessing phase. Actually, when the number of facilities is not too large, one can even require to open 3 or more facilities after having efficiently enumerated all 1- and 2-facility solutions (in our implementation, this latter option is activated only for instances with linear costs and $n \leq 1000$).

Due to the convexity of the $\Phi_j(y)$'s, master problem (12)-(15) is in fact a convex MINLP that can be solved as a MILP by a branch-and-cut approach where linear outer-approximations of constraints (14) are generated on the fly as follows. Let $y^* \in [0, 1]^I$ be a given (possibly noninteger) point with $\sum_{i \in I} y_i^* \geq 2$, and consider a generic customer j . Barring subscript j to ease notation, because of convexity, function $\Phi(y)$ can be underestimated by a supporting hyperplane at y^* , so we can write

$$w \geq \Phi(y) \geq \Phi(y^*) + \sum_{i \in I} \frac{\partial \Phi(y^*)}{\partial y_i} (y_i - y_i^*) \quad (16)$$

Note that (16) can be interpreted as a sensitivity analysis of $\Phi(y)$ in y^* , so it is not surprising that its computation actually requires dual information according to the following scheme.

Each term $\Phi(y)$ is computed by solving a convex *slave subproblem*. For the sake of generality, let this slave be generically written as

$$\Phi(y) = \min\{f(x, y) : g_k(x, y) \leq 0, k = 1, \dots, K, x \in X\} \quad (17)$$

where f and g_1, \dots, g_K are convex and twice differentiable functions and X is a closed convex set.

Given y^* , one can solve efficiently the slave problem (17) for $y = y^*$. Let x^* be the optimal primal solution found, and let $u_k^* \geq 0$ be the optimal dual variables associated with $g_k(x, y) \leq 0$. (In our notation, dual variables correspond to Lagrangian multipliers and are nonnegative even for minimization problems with \leq constraints.)

Using Lagrangian duality and KKT conditions, and assuming constraint qualifications hold, Geoffrion [14] (see also [1]) proved that

$$\frac{\partial \Phi(y^*)}{\partial y_i} = \frac{\partial f(x^*, y^*)}{\partial y_i} + \sum_{k=1}^K u_k^* \frac{\partial g_k(x^*, y^*)}{\partial y_i} \quad (18)$$

An intuitive explanation of this result is as follows. By Lagrangian duality and because of convexity (that plays a fundamental role here), for a given optimal dual vector (u_1^*, \dots, u_K^*) the local behavior of $\Phi(y)$ in y^* is determined by the Lagrangian function in u^* , namely for y sufficiently close to y^* one has

$$\Phi(y) \approx \min\{f(x, y) + \sum_{k=1}^K u_k^* g_k(x, y) : x \in X\} = f(x^*, y^*) + \sum_{k=1}^K u_k^* g_k(x^*, y)$$

hence taking partial derivatives of the the right-hand side leads to (18).

The above considerations lead to the following *Generalized Benders (GB) cuts* [14] to be used within a branch-and-cut MIP solver:

$$w \geq \Phi(y^*) + \sum_{i \in I} \alpha_i^* (y_i - y_i^*)$$

where

$$\alpha_i^* = \frac{\partial f(x^*, y^*)}{\partial y_i} + \sum_{k=1}^K u_k^* \frac{\partial g_k(x^*, y^*)}{\partial y_i}$$

2.1 Benders cuts for linear case

For the linear case, Benders decomposition has been already mentioned in [33] along with the fact that the family of Benders cuts is of polynomial size. However, for the last 30+ years, this result has been forgotten and, to the best of our knowledge and to our surprise, no computational studies were done to asses the practical usefulness of this approach.

For the linear case, a generic slave has the form of the following knapsack problem (in minimization form)

$$\Phi(y^*) = \min \sum_{i \in I} c_i x_i \quad (19)$$

$$\text{s.t. } \sum_{i \in I} x_i = 1 \quad (20)$$

$$x_i \leq y_i^* \quad \forall i \in I \quad (21)$$

$$x_i \geq 0 \quad \forall i \in I \quad (22)$$

For a given $y^* \in [0, 1]^I$ with $\sum_{i \in I} y_i^* \geq 2$, let x^* be an optimal primal solution, β^* be the optimal dual variable for constraint (20), and u_i^* be nonnegative optimal dual variables for constraints (21). The Lagrangian function in (x^*, β^*, u^*) then reads

$$\sum_{i \in I} c_i x_i^* + \beta^* (1 - \sum_{i \in I} x_i^*) + \sum_{i \in I} u_i^* (x_i^* - y_i)$$

hence

$$\frac{\partial \Phi(y^*)}{\partial y_i} = -u_i^*$$

so the Benders cut reads

$$w \geq \Phi(y^*) - \sum_{i \in I} u_i^* (y_i - y_i^*)$$

and only requires the computation of the optimal dual variables $u_i^* \geq 0$ associated with constraints (21).

It is well known from knapsack theory that the following *Dantzig algorithm* produces an optimal primal solution x^* and optimal dual variables $u_i^* \geq 0$ associated with constraints (21) for a given $y^* \in [0, 1]^I$ with $\sum_{i \in I} y_i^* \geq 2$.

Let us assume costs have been sorted to get $c_1 \leq \dots \leq c_n$, and let the index of the *critical item* be defined as the index $k \in I$ such that

$$\sum_{i=1}^{k-1} y_i^* < 1 \leq \sum_{i=1}^k y_i^*$$

Then an optimal primal solution x^* is obtained by setting

$$x_i^* = \begin{cases} y_i^* & \text{for } i < k \\ 1 - \sum_{i=1}^{k-1} y_i^* & \text{for } i = k \\ 0 & \text{for } j > k \end{cases} \quad i \in I$$

As to the dual solution, the optimal dual variable for constraint (20) is $\beta^* = c_k$, while the optimal dual variables for constraint (21) are $u_i^* = c_k - c_i$ for $i < k$, and $u_i^* = 0$ for $i \geq k$.

Optimality comes from the feasibility of the primal and dual solutions for their problems, and from the fact that the primal cost $\sum_{i \in I} c_i x_i^* = \sum_{i=1}^{k-1} c_i y_i^* + c_k (1 - \sum_{i=1}^{k-1} y_i^*)$ is equal to the dual cost $\beta^* - \sum_{i=1}^n u_i^* y_i^* = c_k - \sum_{i=1}^{k-1} (c_k - c_i) y_i^*$. It then follows that

$$\Phi(y^*) = c_k - \sum_{i=1}^{k-1} (c_k - c_i) y_i^*$$

so there are only n distinct Benders cuts for a given customer j , one for each $k \in I$, each of the form

$$w + \sum_{i=1}^{k-1} (c_k - c_i) y_i \geq c_k \quad \forall k \in I$$

(assuming locations have been permuted so as to have $c_1 \leq \dots \leq c_n$).

2.2 Generalized Benders cuts for quadratic case

In this case, for a given y^* , the slave is the convex program (recall that $c_i > 0$ for all $i \in I$).

$$\Phi(y^*) = \min \sum_{i \in I} c_i z_i \tag{23}$$

$$\text{s.t.} \quad \sum_{i \in I} x_i = 1 \tag{24}$$

$$x_i \leq y_i^* \quad \forall i \in I \tag{25}$$

$$x_i^2 \leq z_i y_i^* \quad \forall i \in I \tag{26}$$

$$x_i \geq 0 \quad \forall i \in I \tag{27}$$

$$z_i \geq 0 \quad \forall i \in I \tag{28}$$

To avoid technicalities, we assume $y_i^* > 0$ for all $i \in I$, and refer to Section 3.1 for the general case. Also, we relax $=$ into \geq in (24), and we observe that non-negativity constraints (27) and (28) are redundant and can be relaxed.

Let (x^*, z^*) be an optimal primal solution, and β^* , u_i^* and v_i^* be nonnegative optimal dual variables for constraints (24), (25) and (26), respectively. The Lagrangian function in $(x^*, z^*, \beta^*, u^*, v^*)$ reads

$$\sum_{i \in I} c_i z_i^* + \beta^* (1 - \sum_{i \in I} x_i^*) + \sum_{i \in I} u_i^* (x_i^* - y_i) + \sum_{i \in I} v_i^* (x_i^{*2} - z_i^* y_i)$$

hence

$$\frac{\partial \Phi(y^*)}{\partial y_i} = -u_i^* - v_i^* z_i^*$$

and the GB cut reads

$$w \geq \Phi(y^*) - \sum_{i \in I} (u_i^* + v_i^* z_i^*) (y_i - y_i^*) \tag{29}$$

As in the linear case, one can derive a specialized algorithm for a fast (and numerically accurate) solution of the slave problem, and for the construction of the corresponding GB cut (29).

Observe that, because of the positive costs in (23) and of (26), the optimal z^* satisfies $z_i^* = (x_i^*)^2/y_i^*$ for all $i \in I$. Being y^* fixed, we can just remove those variables and work on the associated problem

$$\Phi(y^*) = \min \sum_{i \in I} \gamma_i x_i^2 \quad (30)$$

$$\text{s.t. } \sum_{i \in I} x_i \geq 1 \quad (31)$$

$$x_i \leq y_i^* \quad \forall i \in I \quad (32)$$

with respect to the new ‘‘perspective costs’’ $\gamma_i = c_i/y_i^* > 0$.

Once the primal solutions x^* and the optimal dual variables u^* associated with (32) have been determined, the optimal primal variables z_i^* for model (23)-(28) are easily computed as

$$z_i^* = (x_i^*)^2/y_i^* \text{ for all } i \in I$$

while the optimal dual variables $v_i^* \geq 0$ corresponding to (26) can be computed as c_i/y_i^* so as to satisfy the first order condition $c_i - v_i^* y_i^* = 0$.

Finally, the most violated GB cut for the given y^* can be computed as in (29).

Implementing an ad-hoc QP algorithm

Fast algorithms for solving simple QP problems with a knapsack-like constraint and box constraints are available in the literature (see, e.g., [21, 36]). To make our results reproducible, we next describe the actual algorithm we used for solving QP problem (30)-(32) and, for the sake of completeness, provide a proof of its correctness.

Let us consider first the following *residual problem* where the x variables have no upper bounds, $R \subseteq I$ is a given *residual set* of facilities, and $r \in [0, 1]$ is a given *residual request*:

$$\min \sum_{i \in R} \gamma_i x_i^2 \quad (33)$$

$$\text{s.t. } \sum_{i \in R} x_i \geq r \quad (34)$$

Let $\beta^* \geq 0$ be the dual multiplier for inequality (34). It is well known [18]—and easy to prove using KKT optimality conditions—that an optimal primal-dual pair (x^*, β^*) can be computed as:

```

Input : A point  $y^* > 0$  with  $\sum_{i \in I} y_i^* \geq 2$  along with the cost vector
           $\gamma > 0$ 
Output: Optimal value  $v_{\text{opt}}$ , optimal primal solution  $x^* \geq 0$ , and
          optimal dual solution  $(u^*, \beta^*) \geq 0$  for the quadratic model
          (30)-(32)

1 /* compute optimal primal solution  $x^*$  along with  $\beta^*$  */
2  $R \leftarrow I$ ;
3  $r \leftarrow 1$ ;
4  $\text{again} \leftarrow \text{TRUE}$ ;
5 while (  $r > 0$  and  $\text{again}$  ) do
6    $\beta^* \leftarrow 2r / \sum_{i \in R} (1/\gamma_i)$ ;
7    $\text{again} \leftarrow \text{FALSE}$ ;
8   foreach  $i \in R$  do
9      $x_i^* \leftarrow \beta^* / (2\gamma_i)$ ;
10    if (  $x_i^* > y_i^*$  ) then
11       $x_i^* \leftarrow y_i^*$ ;
12       $r \leftarrow r - y_i^*$ ;
13       $R \leftarrow R \setminus \{i\}$ ;
14       $\text{again} \leftarrow \text{TRUE}$ ;
15    end
16  end
17 end

18 /* compute optimal value  $v_{\text{opt}}$  and dual solution  $u^*$  */
19 if (  $r \leq 0$  or  $R = \emptyset$  ) then  $\beta^* \leftarrow 0$ ;
20  $v_{\text{opt}} \leftarrow 0$ ;
21 for  $i \in I$  do
22   if  $i \in R$  then  $u_i^* \leftarrow 0$  else  $u_i^* \leftarrow \beta^* - 2\gamma_i y_i^*$ ;
23    $v_{\text{opt}} \leftarrow v_{\text{opt}} + \gamma_i (x_i^*)^2$ ;
24 end

```

Algorithm 1: Our specialized QP solver

$$\beta^* = \frac{2r}{\sum_{i \in R} 1/\gamma_i} \quad x_i^* = \frac{\beta^*}{2\gamma_i} \text{ for all } i \in R \quad (35)$$

Our algorithm to solve the quadratic model with bounded variables (30)-(32) derives the optimal solution by subsequently solving the residual problem (33)-(34) for various values of r and R ; see Algorithm 1.

The primal optimal solution x^* is computed at Steps 1-17 together with β^* . To this end, we first set $R = I$ and $r = 1$ (Steps 2-3) and use (35) to compute the optimal values of x_i^* (and β^*) by disregarding their upper bounds in (32). If it happens that all upper bounds are in fact fulfilled, we are done. Otherwise, for each $x_i^* > y_i^*$ we clip x_i^* to y_i^* (Step 11), update r and R accordingly (Steps 12-13), and repeat. The optimal nonnegative dual variables u_i^* corresponding to (32) are finally defined at Step 22.

To prove the correctness of Algorithm 1, we need the following

Lemma 2.1. *Values β^* computed at each iteration of Step 6 define a strictly monotonically increasing sequence.*

Proof. Let us assume, without loss of generality, that in a given iteration only a single i is removed from R . Let β' and β'' be the value of β^* at Step 6 before and after the removal of i from R arising at Step 13, respectively, and let r' and r'' be the corresponding values of r . To simplify notation, let us define $\rho_i = 1/\gamma_i$ and $d = \sum_{t \in R} (1/\gamma_t)$. We have to show that $\beta'' > \beta'$, where $\beta' = 2r'/d$ and $\beta'' = 2(r' - y_i^*)/(d - \rho_i)$. To this end, observe that the removal of i from R implies $x_i^* = \beta'/(2\gamma_i) > y_i^*$ at Step 10, i.e., $2y_i^* < \beta'\rho_i$ holds. So

$$\beta'' = \frac{2r' - 2y_i^*}{d - \rho_i} > \frac{2r' - \beta'\rho_i}{d - \rho_i} = \beta' \frac{2r'/\beta' - \rho_i}{d - \rho_i} = \beta' \frac{d - \rho_i}{d - \rho_i} = \beta'$$

as claimed. \square

Theorem 2.1. *Algorithm 1 returns an optimal primal solution x^* and an optimal dual solution (u^*, β^*) of problem (30)-(32).*

Proof. The Lagrangian function for problem (30)-(32) is defined as

$$L(x, u, \beta) = \sum_{i \in I} \gamma_i x_i^2 + \beta \left(1 - \sum_{i \in I} x_i\right) + \sum_{i \in I} u_i (x_i - y_i^*)$$

As we have a convex quadratic problem with linear inequalities, all we have to show is that (x^*, u^*, β^*) satisfies the KKT conditions:

$$\nabla_x L(x^*, u^*, \beta^*) = 0 \quad (36)$$

$$(x^*, u^*, \beta^*) \geq 0 \quad (37)$$

$$\sum_{i \in I} x_i^* \geq 1 \quad (38)$$

$$\beta^* \left(1 - \sum_{i \in I} x_i^*\right) = 0 \quad (39)$$

$$u_i^* (x_i^* - y_i^*) = 0 \quad \forall i \in I \quad (40)$$

Let notation R and r refer to the quantities available at the end of the algorithm. For condition (36), we have two cases:

- $i \in R$: in this case, $u_i^* = 0$ (Step 22) and $x_i^* = \beta^* / (2\gamma_i)$ (Step 9), hence $\partial L(x^*, u^*, \beta^*) / \partial x_i = 2\gamma_i x_i^* - \beta^* = 0$;
- $i \notin R$: in this case, $u_i^* = \beta^* - 2\gamma_i y_i^*$ (Step 22) and $x_i^* = y_i^*$ (Step 11), hence $\partial L(x^*, u^*, \beta^*) / \partial x_i = 2\gamma_i x_i^* - \beta^* + u_i^* = 0$.

Conditions (37) are obvious except for u_i^* and $i \notin R$, when $u_i^* = \beta^* - 2\gamma_i y_i^*$ is defined at Step 22. Let β' be the value of β^* computed at Step 6 at the iteration when i is removed from R at Step 13. Then $x_i^* = \beta' / (2\gamma_i) > y_i^*$ at Step 10. Due to Lemma 2.1, $\beta^* \geq \beta'$. By combining this with the rewritten form $\beta' > 2\gamma_i y_i^*$ of the previous inequality, we obtain $\beta^* > 2\gamma_i y_i^*$, from which $u_i^* \geq 0$ follows.

Also obvious is the complementary slackness condition (39), while (40) derives from $u_i^* = 0$ for all $i \in R$ (Step 22) and $x_i^* = y_i^*$ for all $i \notin R$ (Step 11). \square

Thinning out w_j 's variables in the master

In the aim of thinning out irrelevant variables from the master, one can think of replacing all w_j 's by a single continuous variable

$$w_{sum} = \sum_{j \in J} w_j$$

and of aggregating the individual GB cuts (2) for each customer $j \in J$, namely

$$w_j \geq \Phi_j(y^*) + \sum_{i \in I} \alpha_{ij}^* (y_i - y_i^*) \quad (41)$$

into a single *cumulative GB cut*

$$w_{sum} \geq \sum_{j \in J} \Phi_j(y^*) + \sum_{i \in I} \left(\sum_{j \in J} \alpha_{ij}^* \right) (y_i - y_i^*) \quad (42)$$

In this setting, for each y^* of the master one generates (at most) one violated cut, as opposed to the (at most) $|J|$ individual cuts that can be generated by keeping the disaggregated variables w_j 's.

In what follows, the model with the individual w_j 's variables will be called the *fat model*, while the model with the single w_{sum} variable will be referred to as the *slim model*.

The slim model has both pros and cons with respect to the fat one.

From the positive side, one does not really need to know the individual w_j 's values associated with a given y , as the actual solution cost only depends on their sum. By removing the individual w_j 's one can therefore remove potentially-disturbing information. In addition, by using the single w_{sum} variable one expects to generate fewer cuts during the branch-and-cut solution process, hence

the LP relaxations of the master will be smaller in terms of variables and (hopefully) of explicitly-generated GB cuts.

From the negative side, cuts (42) tend to be denser than (41). In addition, we know that each cumulative cut (42) is just implied (à la Farkas) by their individual counterparts (41), and people working with MIP’s have a Pavlov’s unconditioned aversion to aggregated formulations. However, in our setting the quality of the lower bound of the master is theoretically the same as they both are the projection on different w -subspaces of the same formulation, so we are not really losing anything in terms of lower bound if we opt for the cumulative cut.

As it will be discussed in the next section, what matters is in fact the availability of a sound cutting plane scheme that can effectively produce good lower bounds even when a single cut is generated at each separation call. As a matter of fact, at least for quadratic UFL, the slim variant is much more effective than the fat one; see Subsection 3.4 for details.

3 The overall solution framework

Once the GB cut separators for both the linear and quadratic cases have been implemented, one has to design the overall solution framework for the master MILP problem.

A natural choice is to use a state-of-the-art commercial MILP solver—we used IBM ILOG Cplex 12.6 in our implementation. Using a state-of-the-art MILP solver does in fact simplify a lot the implementation, as one relies on a very robust and efficient external framework for the parametric solution of the node LP’s, primal heuristics, generation of general MILP cuts, branching, cut pool handling, etc. So, for a quick shot one only has to plug the GB cut separator in the framework, and see how it works.

This simple approach is in fact rather effective, as it already allowed us to solve previously unsolved UFL instances for both linear and quadratic cases. However, we obtained much better results (in particular, for the quadratic case) by adding some more ingredients to the basic recipe, as outlined below.

3.1 Numerically accurate GB cuts for the quadratic case

In the quadratic case, when the point y^* to separate contains zero entries we face theoretical issues in the definition of the “perspective costs” $\gamma_i = c_i/y_i^*$ and hence of the most-violated GB cut. Although in theory there are elegant ways to cope with the non-differentiability at such points (see [19]), in practice we observed that the GB cut becomes numerically unreliable even when $y_i^* \approx 0$ for some i , hence the risk of producing invalid cuts becomes a real issue.

We have therefore implemented a very simple way to encompass the above difficulty which is based on the fact that, in our separation framework, y^* is in fact a known quantity—while in the compact MILP (5)-(11) it plays the role of a variable. We introduced a (not too small) internal threshold $\varepsilon = 10^{-5}$ to

decide whether a given y^* is close enough to integrality. Then, for each $i \in I$ with $y_i^* < \varepsilon$ we just resist the temptation and do not apply the perspective strengthening of x_{ij}^2 . This means that we replace $x_{ij}^2 \leq z_{ij} y_i$ by $x_{ij}^2 \leq z_{ij}$ in (26). In this way we have $\gamma_i = c_i$ in (30), and the dual variable v_i^* will no longer appear in the GB cut as the new constraint does not depend on y anymore.

Although not strictly required, we also apply the same procedure when $y_i^* > 1 - \varepsilon$, as in this case what we gain from the perspective strengthening is negligible. In a sense, we view the components of y^* that are “almost integer” as good enough, and we do not believe it is worth to penalize them through the (numerically risky) perspective reformulation.

There is a second way to deal with zero entries in y^* , that we call 2ε *trick*: just replace each y_i^* with $y_i^* + 2\varepsilon$ and apply GB separation to this perturbed point—an idea used by many authors, including [5]. This is mathematically correct because the GB cut we generate is always valid—though its violation with respect to original y^* is slightly underestimated. Although it may appear naive, this approach is rather effective and can significantly improve the overall convergence of the cut loop, as shown in the next subsection. In fact, this idea is closely related to the approach proposed by [17] of “perturbing” perspective inequalities by moving y towards a point “sufficiently inside” the perspective cone—the main difference being that we face a much simpler situation where y^* is given, whereas in [17] one needs to treat y as a variable.

3.2 Cut loop stabilization and the curse of Kelley

This is perhaps the most important ingredient in our recipe. In the MIP community, a lot of attention is generally paid to the polyhedral strength of the generated cuts (e.g., the fact of being facet defining), giving for granted that the external cutting plane loop will follow Kelley’s scheme [23]. According to this scheme, at each cut loop iteration one generates one or more cuts that are violated by the current (fractional) solution y^* , adds them to the current relaxation, reoptimizes it and gets a new optimal solution y^* to be cut at the next iteration.

However, the convergence behavior of the overall cut loop heavily depends also on the strategy for generating the next point to cut. As a matter of fact, the famous ellipsoid method has a very good (theoretical) convergence property, but the single valid cut it generates at each iteration can be very weak in polyhedral terms (actually, it does not even need to be violated but just tight at the separation point). This is because the role of the cut is not to let integer points emerge as vertices of the LP relaxation, but just to exclude a large subspace from further considerations. Evidently, one can design a sound cutting plane loop even with very weak/dense cuts, provided that the “Kelley’s curse” is escaped and a different cutting plane scheme is adopted.

In the LP relaxation of our MILP master problem, we have a very simple set of constraints in the y space (essentially, only the 0-1 bounds on the variables), and we have to minimize the convex function $\sum_{j \in J} \Phi_j(y)$. This setting is very close (in fact, identical) to the one arising in the usual Lagrangian dual

minimization (assuming a convex primal problem in maximization form), where “stabilized” approaches such as the bundle method [29] are known to outperform Kelley’s one by a large margin. Thus, the implementation of stabilized cutting plane at the root node (at least) is expected to be of crucial importance, in particular when a single (possibly weak) cut is separated for each point, as it is the case of our slim master model with a single w_{sum} variable.

In our current version, we did not implement a real bundle method, but a simple in-out variant very much in the spirit of [5, 11]. At each cut loop iteration, we have two points: the optimal solution (y^*, w^*) of the current master LP (as in Kelley’s method), and a stabilizing point (\tilde{y}, \tilde{w}) which is initialized to $\tilde{y} = (1, \dots, 1)$ and $\tilde{w} = (\Phi_1(\tilde{y}), \dots, \Phi_m(\tilde{y}))$ (or to $\tilde{w}_{sum} = \sum_{j \in J} \Phi_j(\tilde{y})$ in case the slim model is used).

At each step, we move (\tilde{y}, \tilde{w}) towards (y^*, w^*) by setting

$$(\tilde{y}, \tilde{w}) = 0.5(\tilde{y}, \tilde{w}) + 0.5(y^*, w^*)$$

and then apply our GB separator in the attempt to cut the “intermediate point”

$$\lambda(y^*, w^*) + (1 - \lambda)(\tilde{y}, \tilde{w}) + \delta(1, \dots, 1)$$

Parameters $\lambda \in (0, 1]$ and $\delta \geq 0$ are initially set to 0.2 and 2ε , meaning that we initially try to cut off points very close to the stabilizer. If a violated cut is found, it is statically added to the current LP. After 5 consecutive iterations in which the LP bound does not improve, parameter λ is reset to 1 and the cut loop continues. After 5 more consecutive iterations with no LP bound improvement, parameter δ is reset to 0 (so we are back to Kelley’s scheme). After 5 more consecutive iterations without improvement, the procedure is aborted and all cuts with a positive slack in the final LP are removed. To speedup computation, slack cuts from the current LP are also removed at every 5-th iteration.

According to our computational experience, the above cut loop is very effective for the slim model where GB separation returns (at most) one single cut at each call, in particular for the quadratic case where the difference with respect to the straightforward Kelley’s is striking.

In Figure 1 we plot the behavior of three cut loop strategies (single-thread run) on a sample instance with quadratic costs, namely, the Koerkel-Ghosh [25] instance **gs250a-1** with 250 locations and 250 clients whose optimal value is 12,633.858555. All methods are based on exactly the same GB separator, and only differ on the policy to select the point y^* to cut at each iteration. **Kelley** refers to the standard approach, while **Kelley+** adopts the 2ε trick and slightly perturbs the point to cut before invoking the GB separator. Finally, **inout** refers to our simple scheme using the stabilizer (\tilde{y}, \tilde{w}) as described above. Note that in the horizontal axis we had to use a logarithmic scale due to the dramatically different performance of the three methods. We consider both the fat (top subfigure) and the slim (bottom subfigure) models.

When the fat model is used (top subfigure), **Kelley** has a very poor performance: we stopped it after 15,000 sec.s (still at the root node) with a lower bound of just 1,781.035992, and more than 100,000 cuts generated. **Kelley+**

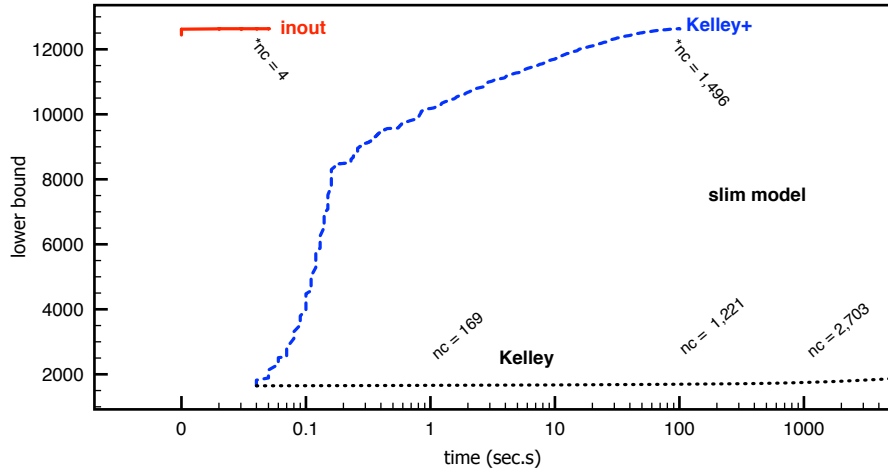
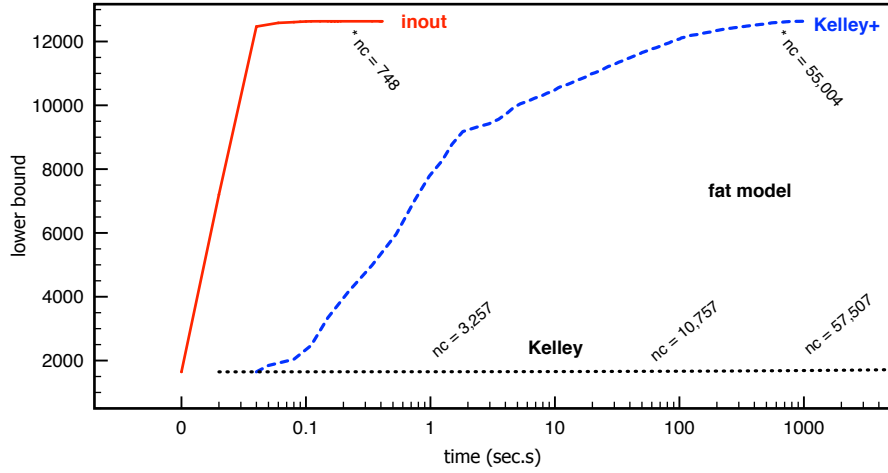


Figure 1: Performance of three cut loop strategies with fat (top subfigure) and slim (bottom subfigure) models on the sample Koerkel-Ghosh [25] instance `gs250a-1` with $n = m = 250$ (quadratic costs). We compare standard Kelley, Kelley+ (i.e., Kelley enhanced by the 2ε trick), and our `inout` scheme. Label `*nc` reports the number of generated cuts at the end of the root node. Time axis given in logarithmic scale.

has a much better performance than `Kelley`, showing the importance of the 2ϵ trick. Its root node takes 948.86 sec.s to generate 55,004 cuts, and produces a very tight bound of 12,633.118973. Enumeration to prove optimality takes 48 branching nodes and is completed at time 978.53. The performance of `inout` is however dramatically better: its root node requires just 0.19 sec.s and produces a lower bound of 12,633.280331 with 748 cuts, while enumeration ends at time 0.41 after 6 nodes.

The performance differences are even more striking for the slim model (bottom subfigure). `Kelley` has again a very poor performance: we stopped it after 500 nodes (8,208.67 sec.s) with a lower bound of 1,875.111861 and 6,037 cuts generated. `Kelley+` has a better performance: its root node takes 98.27 sec.s to generate 1,496 cuts, and produces a bound of 12,632.092567; enumeration to prove optimality takes 44 branching nodes and is completed at time 100.66. The performance of `inout` is so good that its plot is barely visible in the figure: its root node requires 0.04 sec.s and produces a lower bound of 12,633.101453 after adding only 4 cuts, while enumeration ends at time 0.06 after 11 nodes.

3.3 Cut loop along the tree

At each branch-decision node, the MILP framework automatically invokes our GB cut separator (within a so-called `user cut callback`) just before branching, after having solved the current node LP and having added possible violated cuts from its internal cut pool and/or generated by internal procedures. The current LP solution (y^*, w^*) is passed to our GB cut separator, that possibly returns one or more violated cuts.

The violated GB cuts returned by our separation procedure, if any, are added to the internal cut pool and eventually to the current-node LP, which is reoptimized to provide a new solution (y^*, w^*) , and the approach is iterated. Thus, the solver natively implements the Kelley cutting plane scheme, which is not necessarily the best possible option for weak/dense cuts. Implementing a more clever cut loop within Cplex is however not immediate, so we preferred to keep the default Kelley’s scheme within the MILP solver.

To avoid tailing off, we imposed a limit of 20 consecutive cut loop iterations for each node (2,000 for the root node).

Although globally valid, at each node all generated cuts are added as “local cuts” as this allows the solver to remove more cuts from its pool—this turned out to be mainly useful in the case with linear costs, where lots of cuts are generated.

3.4 Slim or fat model?

According to our tests, the slim version with w_{sum} variable is a clear winner for the quadratic case. In our view, this is due to two main reasons. First of all, the stabilized cut loop of the previous subsection performs very well for the quadratic case, and allows us to compute a really tight root-node lower bound with a very small LP with just $n + 1$ variables and very few cuts. Second, the

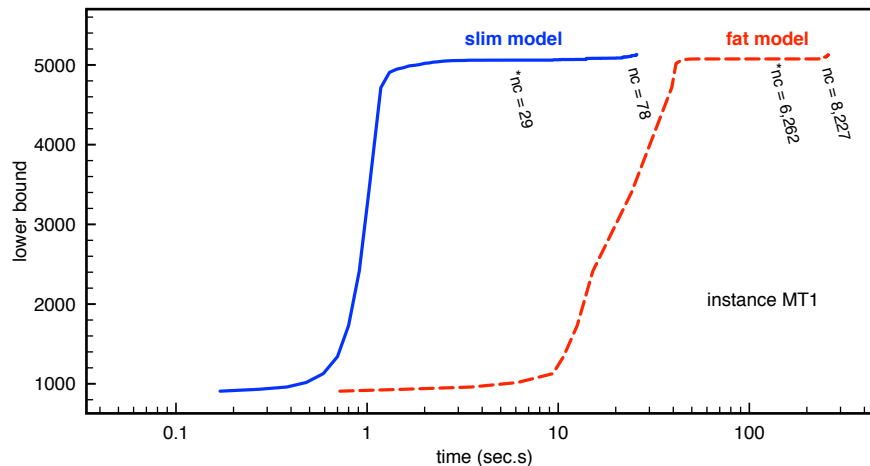


Figure 2: Comparison of the fat and slim models (quadratic costs) on the sample instance MT1 with $n = m = 2000$ when the in-out cut loop is used at the root node. Label *nc reports the number of generated cuts at the end of the root node, while label nc gives the same information at the end of the branch-and-cut algorithm. Time axis is in logarithmic scale.

lower bound is so tight that very few branching nodes are generated, so the fact that the poor Kelley’s cut loop is applied at the non-root nodes is not really an issue.

On the contrary, for the linear case the fat version with the individual w_j ’s variables turns out to be much better. A main reason is that the root-node lower bound is not really tight, so a considerable amount of nodes are enumerated anyway. So Kelley’s curse prevents an effective recomputation of the node bounds when a single cut at the time is generated (w_{sum} version), while it performs reasonably well when several cuts are generated (version with individual w_j ’s). Another reason is that the stabilized cut loop at the root with a single cut at the time (w_{sum} version) is more effective when the convex function to optimize is smooth, which is true only for the quadratic case—while in the linear case Benders’s cuts have a discrete nature.

In Figure 2 we plot the behavior of the fat and slim models on the sample instance MT1 with $n = m = 2000$ and quadratic costs and 4-thread run. At the root node, our in-out cut loop scheme is adopted in both cases. Note that times are reported in logarithmic scale.

3.5 Optimality cuts for integer solutions

Within the MILP branch-and-cut framework, master solutions \tilde{y} with integer \tilde{y}_i ’s can be generated by primal heuristics or when the current-node solution

happens to be integer. In all cases, before updating the incumbent the solution needs to be checked for validity, hence it is passed to a so-called **lazy cut callback** that certifies its validity, or returns one or more valid cuts that prevent the incumbent update. In the latter case, the violated cuts could be obtained by applying our GB separator, that however can return several cuts in case the model with individual w_j 's is used, thus overloading the cut pool with cuts that are likely be active only at the current point \tilde{y} .

Instead, we found it is computationally more efficient to generate a single *optimality cut for integer solutions* within the **lazy cut callback**, namely

$$w_{sum} \geq \left(\sum_{j \in J} \Phi_j(\tilde{y}) \right) \left(1 - \sum_{i \in I: \tilde{y}_i = 0} y_i \right) \quad (43)$$

stating the obvious property that an allocation cost smaller than $\sum_{j \in J} \Phi_j(\tilde{y})$ can only be obtained by opening one or more additional facilities.

As the cut involves the w_{sum} variable, in the model with individual w_j 's we added variable w_{sum} to the master together with its defining equation $w_{sum} = \sum_{j \in J} w_j$.

Of course, GB cuts are still generated within the **user cut callback** for cutting the fractional optimal LP solutions at the various branching nodes.

3.6 Primal heuristics

For the linear case we implemented the following primal heuristics:

- a) at each branching node, within the so-called **heuristic callback**, a simple rounding heuristic is applied. Given the current LP solution y^* , we consider all possible thresholds $\theta \in \{y_i^* : i \in I\}$, in increasing order, round down all y_i^* 's below θ and up all the other y_i^* 's, and evaluate the cost of the integer solution found. By using a parametric technique for cost recomputation (working for the linear case only), this approach just requires $O(\sigma \log \sigma + \sigma m)$ time, where σ is the number of nonzero entries in y^* , so it is very fast.
- b) local branching: right after the cut loop at the root node, we apply the local branching heuristic [9] with neighborhood radius starting from $k = 5$ and then increased to 10, using a small time limit for each call (5000 ticks, corresponding to approximately 5 sec.s on our hardware). The heuristic is aborted when no improved solution is found in the neighborhood of size 10.
- c) proximity search: right after local branching, we apply proximity search heuristic [10] until no improved solution can be found within the small time limit imposed for each call (5000 ticks).

As to the quadratic case, according to our experience the performance of our method is so good that there is not really need to design specific primal heuristics. So for the quadratic case only the rounding heuristic a) above is applied, with a fixed threshold $\theta = \max\{y_i^* : i \in I\} - 0.2$.

3.7 Numerical tolerances and cut validity

In our implementation, we were very conservative and used very small tolerances for numerical and integrality tests. To be specific, we set Cplex’s integrality tolerance `CPX_PARAM_EPINT` to 0, and the optimality/violation tolerances `CPX_PARAM_EPGAP` and `CPX_PARAM_EPRHS` to 10^{-9} . Our own internal numerical tolerance was set to 10^{-9} as well.

To assert the validity of the GB cuts we generate in our code, we implemented the following check in the spirit of Margot’s proposal [34]. At the end of the run (or at the time limit), we fix all (binary and continuous) variables to their value in the incumbent solution, and enter a cut loop where (1) we apply our GB cut separator to a point y^* obtained from the incumbent by a small random perturbation (applied to 80% random entries) ranging from 10^{-9} to 10^{-1} , (2) we statically add the generated cuts to the current MILP, and (3) we optimize the MILP to verify its feasibility. If the current MILP becomes infeasible, we write down the current model in a file for further analysis and report a failure, otherwise we repeat for 10,000 times. The above check was applied extensively during the development of our code, and helped up in detecting and correcting some tolerance issues present in the earlier versions. Needless to say, our final code passed the test for all instances we tried.

4 Computational results

In this section we report on our computational experience on a subset of most difficult instances for UFL. As to qUFL, we consider instances used in the previous literature, extended by a family of much larger instances that cannot be approached by existing methods (mainly due to the fact that the underlying models would consist of millions of variables and constraints). The computational study is conducted on a cluster of identical machines each consisting of an Intel Xeon E3-1220V2 CPU running at 3.10 GHz, with 16GB of RAM each. Reported times reported are wall-clock seconds and refer to 4-thread runs.

We report computational results with basic parameter settings, without tuning our code with respect to proximity search and/or local branching parameters that could theoretically further improve the obtained results.

4.1 Benchmark instances

The set of UFL benchmark instances used in this paper stems from the UFLLIB [22], which is a well-established library of instances for capacitated and uncapacitated facility location problems. The library is a diverse collection of bipartite graphs, some of them being just small and easily solvable cases. In our study on UFL, we focus on a subset instances from UFLLIB representing the most challenging ones even for the most recent state-of-the-art approaches, like the ones proposed by [31, 37, 4]. These are randomly generated instances M^* (proposed in [26]) and KG (proposed in [15, 25]). Instances M^* are of size 100×100 up to 2000×2000 , whereas KG instances can be divided into three groups, with

$n = m \in \{250, 500, 750\}$. Within each KG group, there are two classes, symmetric and asymmetric ones, denoted by gs^* and ga^* , respectively. Additionally, each class contains three subclasses, “a”, “b” and “c”, representing different cost settings: in “a”, allocation costs are an order of magnitude higher than the facility opening costs; in “b”, these costs are of the same order; and in “c”, facility opening costs are an order of magnitude higher than the allocation costs. As we will report below, these differences in the costs structure significantly influence their computationally difficulty.

For testing the impact of Benders decomposition to the separable quadratic case, we consider two families of benchmark instances: (1) UFLLIB instances mentioned above plus the instances from the ORLIB (with original allocation costs equal to 0 being replaced by 10^{-5}), and (2) randomly generated instances used in previous computational studies in [6, 18, 19]. The latter instances are randomly generated graphs with potential facility and customer locations being placed uniformly at random within a unit square, with allocation costs calculated as the Euclidean distance multiplied by a factor of 50, and facility opening costs generated uniformly in [1, 100]. Tables shown in this section and in the Appendix report the total running time in wall-clock seconds ($t[s]$), the time needed to solve the LP-relaxation at the root node ($t_{root}[s]$), the total number of branch and bound nodes needed to prove the optimality (*nodes*) and the percentage gap at the root node ($g_r[\%]$) and the bound at the root node (*rootbound*).

4.2 Linear costs

For solving UFL, we opted for the branch-up-first setting (Cplex’s parameter `CPX_PARAM_BRDIR` set to 1), as this tends to produce branching trees with fewer open nodes—hence reducing the overhead incurred when writing node files.

Table 1 summarizes results for the previously unsolved UFL instances for which we were able to prove optimality. In total, optimal solutions are provided for seven previously unsolved instances.

Two of them belong to the benchmark set of 18 instances proposed by Barahona-Chudak [2]. The semi-Lagrangian approach by [4] solved 16 of them, whereas the approach by [37] managed to solve only 8 of these 16. We provide the optimal values for the remaining two unsolved instances, namely, 2500-10 and 3000-100, of size 2500×2500 and 3000×3000 , respectively.

Concerning the 30 instances of size 250×250 proposed by Koerkel-Gosh [25], 27 were solved by [37], and an additional one ($ga250a-1$) was solved by [4]. We now provide optimal values for the two previously unsolved instances, namely, $ga250a-3$ and $ga250a-5$. Among the KG instances of size 500×500 , optimal values were known only for 7 (out of 10) instances of the subclass $g*500c*$. We were able to prove the optimality for the missing 3 instances from this subclass, namely $ga500c-5$ and $gs500c-3$ and $gs500c-5$.

The 50 KG instances of subclasses $g*500a$, $g*500b$, $g*750*$ still remain out of reach for existing exact methods. However, we managed to improve the best known upper bounds for 22 of these instances. To this end, we slightly

inst.	bestknown	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
ga250a-3	257985	257953	493.49	257554.773407	12.77	0.15	200184
ga250a-5	258225	258190	585.93	257790.245068	9.65	0.15	229446
ga500c-5	621313	621313	9226.86	601500.282332	12.31	3.19	195191
gs500c-3	621204	621204	11448.19	601980.526816	13.44	3.09	194657
gs500c-5	623180	623180	26828.91	603115.401650	14.20	3.22	270147
2500-10	3101800	3099907	824.76	3097480.189279	104.67	0.08	1362
3000-100	1602335	1602154	225.25	1601733.816607	82.67	0.03	441

Table 1: Previously unsolved UFL instances solved to optimality using our approach (linear costs).

modified the initial local-branching heuristic described in Subsection 3.6 by removing the 5000-tick timelimit for each subMIP solution. We implemented two local branching variants: in variant (A), the neighborhood radius is 5 or 10, while in variant (B) the radius is 2, 4, 6, 8 or 10. For both versions, if no improved solution can be found in the neighborhood of size 10, we randomize the incumbent and start again from the smallest radius. To better exploit the 4-core architecture of our PC’s, we concurrently ran 4 times our algorithm in 1-thread mode, with 4 different input random seeds, and collected the best solution found. Our results are summarized in Table 2, and refer to a timelimit of 600 and 3600 wall-clock sec.s for each instance, respectively. Column *bestknown* gives the previous best upper bound from the literature, while *UB* is the heuristic value we computed within the given time limit. The remaining columns report a 0/1 flag saying whether we strictly improved (*win*) or matched (*tie*) the previous best known solution. Finally, columns under the *Best* header refer to the best of the two 3600-sec.s runs with variants (A) and (B). According to the table, 600 sec.s are enough for variant (A) to strictly improve 14 (and match 13) best-known solutions, while in 3600 sec.s we strictly improved 19 (and matched 21) solutions. Running variant (B) for 3600 more sec.s slightly improved our results, and leads to 22 strictly improved (and 22 matched) solutions. For the 6 instances (out of 50) where our *UB* is worse than *bestknown*, the gap between these two values is always below 0.04%.

4.3 Separable quadratic costs

For testing qUFL, we first report results obtained on a set of smaller randomly generated instances created according to the procedure described in [18]. Table 3 reports the average values over 10 instances generated with the fixed number of facilities and customers (shown in columns n and m , respectively). We compare our slim and fat models against the perspective reformulation proposed in [19]. The latter model was solved by Cplex after setting parameter `CPX_PARAM_MIQCPSTRAT` to 1, meaning that a QCP relaxation is solved at each node—this setting turned out to be much better than the default one where cone cuts are generated.

The obtained results clearly demonstrate the power of Benders decomposition for quadratic separable objective functions. Recall that the main bottleneck of the perspective reformulation of [19] is its $O(n \cdot m)$ number of variables and constraints. Table 3 shows that both our formulations scale extremely well with the increasing size of the input data. All problem instances with up to 250 facilities and 250 customers shown in this table could be solved within a fraction of a second. Compared to the performance of the perspective reformulation, our models allow for computational speedups of up to four orders of magnitude. For $n, m \geq 200$, the perspective reformulation already hits the memory limit, and it is even impossible to solve the QP/QCP relaxation at the root node.

In our next experiment, we decided to push our decomposition approaches to their limits, and for that purpose we created a set of much larger instances following the same graph generation procedure used for the instances shown in Table 3. To this end, we considered input graphs with $n \in \{500, 1000, 2000\}$ and $m \in \{500, 1000, 5000, 10000\}$. Table 4 shows the comparison of the performance of our slim and fat models. We notice that fat model is outperformed by slim model, and that the difference in the running times increases with the increasing number of customers. For example, the fat model is only a few times slower than slim for $m \leq n$, but for the larger values of m , slim model is significantly faster than its fat counterpart. Solving even the largest instances of this group (with 2000 facilities and 10000 customers) using our slim approach requires only about 5 minutes on average. The quality of the root gap (that is computed by our in-out algorithm followed by the usual root-node processing) is quite remarkable—the average root gap is consistently below 0.04%. The small difference in the quality of the LP-relaxation gaps between the two models can be explained by tailing off. Comparing the time required to solve the LP-relaxation at the root node, we observe that our both models spent most of their computing time at the root node (except for the largest instances). The required number of branch-and-bound nodes remains quite moderate for all instances, except for the largest ones, where its average value exceeds 10000.

To our big surprise, even the most difficult UFLLIB instances (e.g., those from KG of size 500×500 and 750×750), are easily solvable as qUFL's by our slim model. More precisely, all of M and KG instances can be solved to optimality in about half a minute or much less. The most difficult instance appears to be MT1 (of size 2000×2000) for which our slim model requires just 81.99 seconds. As a comparison, the same instance requires 4889.77 sec.s and 10675 branching nodes when given on input to our UFL code (linear case). This behavior is of course explained by the fact that, for a same input, the lower bounds are typically much tighter for qUFL than for (linear) UFL, so much fewer branching nodes are required.

A summary of the obtained results on KG instances (containing average values for each subclass of five instances) is shown in Table 5. More detailed results are reported in the Appendix.

instance	bestknown	Local Branching (A)						Local Branching (B)			Best		
		600 seconds			3600 seconds			3600 seconds			7200 seconds		
		<i>UB</i>	win	tie	<i>UB</i>	win	tie	<i>UB</i>	win	tie	<i>UB</i>	win	tie
ga500a-1	511422	511401	1	0	511383	1	0	511388	1	0	511383	1	0
ga500a-2	511333	511288	1	0	511255	1	0	511255	1	0	511255	1	0
ga500a-3	510817	510810	1	0	510810	1	0	510810	1	0	510810	1	0
ga500a-4	511047	511008	1	0	511008	1	0	511008	1	0	511008	1	0
ga500a-5	511258	511239	1	0	511239	1	0	511239	1	0	511239	1	0
ga500b-1	538060	538656	0	0	538060	0	1	538060	0	1	538060	0	1
ga500b-2	537850	537850	0	1	537850	0	1	537850	0	1	537850	0	1
ga500b-3	538077	538144	0	0	537924	1	0	537924	1	0	537924	1	0
ga500b-4	537925	538038	0	0	537925	0	1	537925	0	1	537925	0	1
ga500b-5	537482	537642	0	0	537482	0	1	537482	0	1	537482	0	1
gs500a-1	511229	511201	1	0	511188	1	0	511188	1	0	511188	1	0
gs500a-2	511179	511179	0	1	511179	0	1	511179	0	1	511179	0	1
gs500a-3	511120	511129	0	0	511112	1	0	511112	1	0	511112	1	0
gs500a-4	511137	511137	0	1	511137	0	1	511137	0	1	511137	0	1
gs500a-5	511293	511293	0	1	511293	0	1	511293	0	1	511293	0	1
gs500b-1	537931	537941	0	0	537931	0	1	537931	0	1	537931	0	1
gs500b-2	537763	537823	0	0	537763	0	1	537763	0	1	537763	0	1
gs500b-3	537874	538095	0	0	537926	0	0	537854	1	0	537854	1	0
gs500b-4	537742	537779	0	0	537742	0	1	537779	0	0	537742	0	1
gs500b-5	538270	538270	0	1	538270	0	1	538270	0	1	538270	0	1
ga750a-1	763576	763537	1	0	763537	1	0	763528	1	0	763528	1	0
ga750a-2	763674	763679	0	0	763653	1	0	763674	0	1	763653	1	0
ga750a-3	763765	763748	1	0	763697	1	0	763699	1	0	763697	1	0
ga750a-4	764033	764043	0	0	763945	1	0	763976	1	0	763945	1	0
ga750a-5	763905	763857	1	0	763794	1	0	763786	1	0	763786	1	0
ga750b-1	796480	796506	0	0	796454	1	0	796454	1	0	796454	1	0
ga750b-2	796056	796003	1	0	795963	1	0	795963	1	0	795963	1	0
ga750b-3	796130	796439	0	0	796384	0	0	796359	0	0	796359	0	0
ga750b-4	797080	797013	1	0	797013	1	0	797013	1	0	797013	1	0
ga750b-5	796387	796549	0	0	796549	0	0	796549	0	0	796549	0	0
ga750c-1	902026	902026	0	1	902026	0	1	902026	0	1	902026	0	1
ga750c-2	899651	899732	0	0	899651	0	1	899732	0	0	899651	0	1
ga750c-3	900010	900019	0	0	900019	0	0	900019	0	0	900019	0	0
ga750c-4	900044	900044	0	1	900044	0	1	900044	0	1	900044	0	1
ga750c-5	899235	899235	0	1	899235	0	1	899235	0	1	899235	0	1
gs750a-1	763671	763683	0	0	763683	0	0	763671	0	1	763671	0	1
gs750a-2	763548	763590	0	0	763552	0	0	763548	0	1	763548	0	1
gs750a-3	763764	763759	1	0	763748	1	0	763727	1	0	763727	1	0
gs750a-4	763887	763942	0	0	763932	0	0	763922	0	0	763922	0	0
gs750a-5	763616	763614	1	0	763614	1	0	763616	0	1	763614	1	0
gs750b-1	797026	797688	0	0	797347	0	0	797329	0	0	797329	0	0
gs750b-2	796170	796498	0	0	796170	0	1	796170	0	1	796170	0	1
gs750b-3	796589	796589	0	1	796589	0	1	796589	0	1	796589	0	1
gs750b-4	796734	797020	0	0	797020	0	0	797087	0	0	797020	0	0
gs750b-5	796365	796365	0	1	796365	0	1	796365	0	1	796365	0	1
gs750c-1	900454	900454	0	1	900454	0	1	900363	1	0	900363	1	0
gs750c-2	897886	897886	0	1	897886	0	1	897886	0	1	897886	0	1
gs750c-3	901714	901947	0	0	901786	0	0	901656	1	0	901656	1	0
gs750c-4	901339	901239	1	0	901239	1	0	901239	1	0	901239	1	0
gs750c-5	900216	900216	0	1	900216	0	1	900216	0	1	900216	0	1
sum			14	13		19	21		20	22		22	22

Table 2: Local branching for UFL on the 50 unsolved KG instances (linear costs). We strictly improved 22 (and matched 22 more) best known values from the literature.

n	m	Our slim model				Our fat model				Perspective reformulation [19]			
		$t[s]$	$g_r[\%]$	$t_{root}[s]$	nodes	$t[s]$	$g_r[\%]$	$t_{root}[s]$	nodes	$t[s]$	$g_r[\%]$	$t_{root}[s]$	nodes
10	30	0.01	0.21	0.01	0.0	0.01	0.26	0.00	0.0	1.97	0.00	1.79	3.6
10	50	0.01	0.17	0.01	0.0	0.03	0.17	0.02	0.0	3.59	0.00	3.25	4.4
10	100	0.03	0.30	0.02	2.5	0.09	0.22	0.05	1.8	7.23	0.00	6.21	6.2
10	200	0.03	0.25	0.02	2.0	0.23	0.22	0.13	3.6	16.50	0.00	14.23	6.1
20	30	0.03	0.49	0.01	2.9	0.03	0.42	0.02	2.8	5.20	0.16	4.38	6.1
20	50	0.03	0.39	0.02	3.1	0.05	0.32	0.03	3.0	8.21	0.04	6.98	6.1
20	100	0.04	0.31	0.02	4.5	0.12	0.28	0.06	4.6	19.96	0.19	16.00	8.2
20	200	0.05	0.18	0.03	5.0	0.21	0.15	0.11	5.6	32.91	0.18	23.12	8.7
30	30	0.03	0.38	0.01	1.5	0.03	0.27	0.01	1.0	11.12	0.01	8.75	7.1
30	50	0.03	0.29	0.02	1.4	0.04	0.17	0.03	1.0	17.45	0.00	15.41	3.6
30	100	0.05	0.21	0.03	2.5	0.10	0.22	0.06	2.7	27.51	0.18	19.66	7.1
30	200	0.07	0.25	0.05	4.8	0.26	0.23	0.15	5.6	69.49	0.21	39.58	9.8
40	30	0.03	0.38	0.02	1.8	0.03	0.39	0.01	1.9	17.58	0.00	14.24	5.8
40	50	0.04	0.24	0.02	2.9	0.05	0.22	0.03	3.2	25.52	0.01	20.33	4.9
40	100	0.05	0.22	0.03	3.8	0.12	0.16	0.06	3.8	53.79	0.19	35.76	9.2
40	200	0.09	0.14	0.06	6.1	0.27	0.14	0.12	5.6	104.94	0.16	56.04	10.1
50	50	0.04	0.14	0.03	1.6	0.05	0.13	0.03	1.6	41.08	0.03	33.04	4.4
50	100	0.06	0.13	0.04	2.6	0.10	0.11	0.06	2.7	88.09	0.15	50.02	8.4
50	200	0.11	0.13	0.08	6.7	0.29	0.12	0.16	7.5	159.82	0.14	71.45	11.0
60	50	0.05	0.29	0.03	3.1	0.06	0.27	0.03	2.6	56.64	0.26	37.71	6.4
60	100	0.06	0.12	0.04	1.6	0.10	0.10	0.06	1.5	99.44	0.16	54.41	6.6
60	200	0.12	0.11	0.09	4.6	0.25	0.11	0.15	5.2	280.67	0.14	113.92	15.4
70	30	0.05	0.28	0.03	2.8	0.04	0.26	0.03	2.4	47.65	0.37	25.32	10.5
70	50	0.06	0.23	0.04	3.1	0.06	0.21	0.03	2.7	79.52	0.25	49.09	6.7
70	100	0.09	0.23	0.07	4.3	0.14	0.20	0.09	4.2	186.35	0.27	81.57	10.1
70	200	0.09	0.04	0.08	0.8	0.19	0.03	0.13	1.4	293.19	0.04	176.88	5.9
80	30	0.05	0.22	0.03	1.7	0.05	0.16	0.03	1.1	59.25	0.39	36.57	6.9
80	50	0.08	0.36	0.05	5.7	0.07	0.34	0.04	5.5	108.71	0.57	42.71	13.5
80	100	0.10	0.21	0.08	5.9	0.14	0.21	0.08	5.3	245.21	0.28	104.64	10.0
80	200	0.14	0.13	0.11	5.2	0.27	0.14	0.16	6.4	462.30	2.06	160.77	12.0
100	100	0.23	0.21	0.19	6.6	0.16	0.20	0.10	6.0	482.12	8.63	174.57	15.3
150	150	0.24	0.17	0.19	7.8	0.32	0.16	0.20	9.0	2012.07	3.15	648.82	14.1
200	200	0.33	0.06	0.28	6.7	0.45	0.06	0.32	4.1	—	—	—	—
250	250	0.46	0.05	0.42	4.3	0.71	0.04	0.60	4.1	—	—	—	—

Table 3: Comparing our slim and fat models with the perspective reformulation [19], on a set of randomly generated qUFL instances proposed in [18, 19]. Perspective reformulation hits memory limit for $n, m \geq 200$.

n	m	Our slim model				Our fat model			
		$t[s]$	$g_r[\%]$	$t_{root}[s]$	nodes	$t[s]$	$g_r[\%]$	$t_{root}[s]$	nodes
500	500	1.39	0.03	1.31	16.2	3.30	0.03	2.82	9.5
500	1000	3.02	0.03	2.75	54.7	8.90	0.03	7.81	20.8
500	5000	11.59	0.01	10.41	87.2	132.89	0.02	127.27	32.4
500	10000	36.98	0.01	22.09	558.2	673.93	0.02	646.97	106.5
1000	500	3.80	0.04	3.32	76.0	4.60	0.04	3.86	26.1
1000	1000	5.78	0.03	5.25	65.3	15.18	0.03	13.74	28.2
1000	5000	20.70	0.01	19.32	44.3	193.76	0.02	181.87	180.3
1000	10000	64.01	0.01	34.74	603.0	799.02	0.02	748.56	399.8
2000	500	6.73	0.03	6.10	66.7	8.95	0.03	7.83	29.8
2000	1000	14.86	0.02	12.72	194.4	35.41	0.02	32.65	65.9
2000	5000	115.09	0.01	42.07	1649.0	405.85	0.02	361.69	629.3
2000	10000	309.36	0.01	76.88	10735.8	2646.69	0.03	1246.60	13114.0

Table 4: Comparing the performance of slim versus fat model on a larger set of benchmark instances for qUFL generated as in [18, 19].

5 Conclusions

The Uncapacitated Facility Location (UFL) problem is one of the most famous and studied Operations Research problems. This problem can easily be formulated as a MILP, whose size is however exceedingly large for many practical applications, making the direct use of a MILP solver rather ineffective (or even impossible). As a matter of fact, the most powerful technique at present for solving large scale UFL instances is Lagrangian relaxation, that allows one to quickly compute lower bounds that are close enough to the LP relaxation ones. The implementation of a sound Lagrangian relaxation method is however far from trivial, and a number of sophisticated ideas need to be implemented to get satisfactory results. This was the main motivation for us to consider a simpler approach built on top of an off-the-shelf MILP solver.

We therefore decided to investigate a MILP approach based on Benders decomposition, a technique that can be considered folklore but apparently not used in recent computational studies for UFL. We also addressed a nonlinear version of the problem, namely, separable convex quadratic UFL (qUFL) that has been the subject of intensive computational studies in the recent years.

From the methodological point of view, our approach uses generalized Benders cuts for convex problems, and embeds them in a branch-and-cut scheme. A number of important features are introduced, that are instrumental for the practical effectiveness of the overall approach. In particular, we discuss how to speedup the root-node cut loop through simple stabilization techniques that make it orders of magnitude faster than standard cutting plane loops.

Using our approach, we were able to solve to proven optimality 7 previously unsolved benchmark instances for UFL, and to improve the best-known heuristic value for 22 additional instances. These instances were out of reach for previous

group	$t[s]$	$g_r[\%]$	$t_{root}[s]$	nodes
ga250a	0.40	0.00	0.39	4.4
ga250b	0.28	0.03	0.22	71.2
ga250c	0.40	0.03	0.36	39.4
gs250a	0.21	0.00	0.20	3.4
gs250b	0.27	0.02	0.21	80.6
gs250c	0.46	0.03	0.42	21.6
ga500a	0.76	0.00	0.73	3.0
ga500b	2.12	0.04	1.95	58.0
ga500c	19.46	0.16	1.49	49911.6
gs500a	0.81	0.00	0.77	12.4
gs500b	2.47	0.03	2.31	72.8
gs500c	15.05	0.14	1.26	12721.6
ga750a	2.03	0.00	1.62	107.4
ga750b	2.08	0.01	1.82	65.2
ga750c	35.79	0.08	2.41	64338.0
gs750a	1.94	0.00	1.65	53.2
gs750b	3.24	0.01	1.82	414.0
gs750c	26.94	0.07	2.98	16837.0

Table 5: All KG instances for qUFL are solved to optimality by our slim model (quadratic costs). Each row shows average values over 5 instances per subclass.

MILP approaches, as the underlying models would involve tens of millions of variables and constraints, and they turned out to be too hard even for the best Lagrangian methods from the literature.

Even more interesting results are reported for qUFL: compared to previous methods, our approach enables speedups of 4 orders of magnitude or more, and allowed us to tackle much larger qUFL instances than any previous approach.

The potential impact of our work is twofold:

- On the one hand, we believe our results will motivate researchers to rethink/reinvent decomposition approaches for many optimization problems involving location, allocation and network design decisions. “Thinning out” can be done by reformulating “location subproblems” through a linear number of constraints and variables to model allocation decisions. Problems that can directly benefit from such reformulations and elimination of variables through Benders cuts are: connected facility location [16], the ring-star problem [27], the median cycle problem, or the traveling purchaser problem [28], to mention only a few.
- On the other hand, our specially designed stabilization cutting-plane schemes are of tremendous importance not only for MILPs, but also for the emerging area of convex MINLPs. Exploiting this scheme together with gener-

alized Benders decomposition and/or perspective reformulations may lead to the next boost of performance of convex MINLP solvers.

Although very closely related to UFL, application of our methods to capacitated facility location is not immediate, and will be subject of future research. We will also focus on other non-trivial MINLPs that could benefit from generalized Benders cuts, including nonlinear Stochastic Programming.

Acknowledgments

This research of the first author was supported by the University of Padova (Progetto di Ateneo “Exploiting randomness in Mixed Integer Linear Programming”), and by MiUR, Italy (PRIN project “Mixed-Integer Nonlinear Optimization: Approaches and Applications”). The work of the last author was supported by the Austrian Research Fund (FWF, Project P 26755-N19) and STSM Grant from COST Action TD1207. The authors would like to acknowledge these supports.

References

- [1] E. Balas. A duality theorem and an algorithm for (mixed-) integer nonlinear programming. *Linear Algebra and its Applications*, 4(4):341–352, 10 1971.
- [2] F. Barahona and F. Chudak. Solving large scale uncapacitated facility location problems. In P. Pardalos, editor, *Approximation and Complexity in Numerical Optimization*, pages 48–62. 2000.
- [3] S. Basu, M. Sharma, and P. Ghosh. Metaheuristic applications on discrete facility location problems: a survey. *OPSEARCH*, pages 1–32, 2014.
- [4] C. Beltran-Royo, J.-P. Vial, and A. Alonso-Ayuso. Semi-lagrangian relaxation applied to the uncapacitated facility location problem. *Computational Optimization and Applications*, 51(1):387–409, 2012.
- [5] W. Ben-Ameur and J. Neto. Acceleration of cutting-plane and column generation algorithms: Applications to network design. *Networks*, 49(1):3–17, 2007.
- [6] P. Bonami, M. Kilinc, and J. Linderoth. Algorithms and software for convex mixed integer nonlinear programs. *Mixed Integer Nonlinear Programming*, 154:1–39, 2012.
- [7] C. D’Ambrosio, J. Lee, and A. Wächter. A global-optimization algorithm for mixed-integer nonlinear programs having separable non-convexity. In A. Fiat and P. Sanders, editors, *Algorithms - ESA 2009*, volume 5757 of *Lecture Notes in Computer Science*, pages 107–118. Springer Berlin Heidelberg, 2009.

- [8] M. Fischetti, M. Leitner, I. Ljubić, M. Luipersbeck, M. Monaci, M. Resch, D. Salvagnin, and M. Sinnl. Thinning out Steiner trees: A node-based model for uniform edge costs. *Mathematical Programming Computations*, 2015. Special Issue: 11th DIMACS Implementation Challenge on Steiner Tree Problems, submitted.
- [9] M. Fischetti and A. Lodi. Local branching. *Mathematical Programming*, 98(1-3):23–47, 2003.
- [10] M. Fischetti and M. Monaci. Proximity search for 0-1 mixed-integer convex programming. *Journal of Heuristics*, 20(6):709–731, 2014.
- [11] M. Fischetti and D. Salvagnin. An in-out approach to disjunctive optimization. In A. Lodi, M. Milano, and P. Toth, editors, *Integration of AI and OR Techniques in Constraint Programming for Combinatorial Optimization Problems*, volume 6140 of *Lecture Notes in Computer Science*, pages 136–140. Springer Berlin Heidelberg, 2010.
- [12] A. Frangioni and C. Gentile. Perspective cuts for a class of convex 0-1 mixed integer programs. *Mathematical Programming*, 106(2):225–236, 2006.
- [13] H. A. Friberg. CBLIB 2014: A benchmark library for conic mixed-integer and continuous optimization. *Optimization Online*, March 2014.
- [14] A. Geoffrion. Generalized Benders Decomposition. *Journal of Optimization Theory and Applications*, 10:237–260, 1972.
- [15] D. Ghosh. Neighborhood search heuristics for the uncapacitated facility location problem. *European Journal of Operational Research*, 150(1):150–162, 2003.
- [16] S. Gollowitzer and I. Ljubić. MIP models for connected facility location: A theoretical and computational study. *Computers & Operations Research*, 38(2):435–449, 2011.
- [17] I. E. Grossmann and S. Lee. Generalized convex disjunctive programming: Nonlinear convex hull relaxation. *Computational Optimization and Applications*, 26(1):83–100, 2003.
- [18] O. Günlük, J. Lee, and R. Weismantel. MINLP strenghtening for separable convex quadratic transportation-cost UFL. Technical Report RC24213 (W0703-042), IBM Research Division, 2007.
- [19] O. Günlük and J. Linderoth. Perspective reformulation and applications. In J. Lee and S. Leyffer, editors, *Mixed Integer Nonlinear Programming*, pages 61–92. Springer, 2012.
- [20] H. Hijazi, P. Bonami, and A. Ouorou. An outer-inner approximation for separable mixed-integer nonlinear programs. *INFORMS Journal on Computing*, 26(1):31–44, 2014.

- [21] D. S. Hochbaum and S.-P. Hong. About strongly polynomial time algorithms for quadratic optimization over submodular constraints. *Mathematical Programming*, 69(1-3):269–309, 1995.
- [22] M. Hofer. Uflib, 2006. <http://resources.mpi-inf.mpg.de/departments/d1/projects/benchmarks/UflLib/>.
- [23] J. J. E. Kelley. The cutting-plane method for solving convex programs. *Journal of the Society for Industrial and Applied Mathematics*, 8(4):703–712, 1960.
- [24] A. Klose and A. Drexl. Facility location models for distribution system design. *European Journal of Operational Research*, 162(1):4–29, 2005.
- [25] M. Körkel. On the exact solution of large-scale simple plant location problems. *European Journal of Operational Research*, 39(2):157–173, 1989.
- [26] J. Kratica, D. Tošić, V. Filipović, and I. Ljubić. Solving the simple plant location problem by genetic algorithm. *RAIRO-Operations Research*, 35(01):127–142, 2001.
- [27] M. Labbé, G. Laporte, I. R. Martín, and J. J. S. González. The ring star problem: Polyhedral analysis and exact algorithm. *Networks*, 43(3):177–189, 2004.
- [28] G. Laporte, J. R. Ledesma, and J. S. González. A branch-and-cut algorithm for the undirected traveling purchaser problem. *Operations Research*, 51(66):940–951, 2003.
- [29] C. Lemaréchal, A. Nemirovskii, and Y. Nesterov. New variants of bundle methods. *Mathematical Programming*, 69(1-3):111–147, 1995.
- [30] A. Letchford and S. Miller. Fast bounding procedures for large instances of the simple plant location problem. *Computers & Operations Research*, 39(5):985–990, 2012.
- [31] A. Letchford and S. Miller. An aggressive reduction scheme for the simple plant location problem. *European Journal of Operational Research*, 234(3):674–682, 2014.
- [32] S. Li. A 1.488 approximation algorithm for the uncapacitated facility location problem. *Information and Computation*, 222(0):45–58, 2013. 38th International Colloquium on Automata, Languages and Programming (ICALP 2011).
- [33] T. L. Magnanti and R. T. Wong. Accelerating Benders decomposition: Algorithmic enhancement and model selection criteria. *Operations Research*, 29(3):464–484, 1981.
- [34] F. Margot. Testing cut generators for mixed-integer linear programming. *Mathematical Programming Computation*, 1(1):69–95, 2009.

- [35] M. Melo, S. Nickel, and F. S. da Gama. Facility location and supply chain management – a review. *European Journal of Operational Research*, 196(2):401–412, 2009.
- [36] M. Patriksson and C. Strömberg. Algorithms for the continuous nonlinear resource allocation problem—new implementations and numerical studies. *European Journal of Operational Research*, 2015.
- [37] M. Posta, J. A. Ferland, and P. Michelon. An exact cooperative method for the uncapacitated facility location problem. *Mathematical Programming Computation*, 6(3):199–231, 2014.
- [38] V. Verter. Uncapacitated and capacitated facility location problems. In H. A. Eiselt and V. Marianov, editors, *Foundations of Location Analysis*, volume 155 of *International Series in Operations Research & Management Science*, pages 25–37. Springer, 2011.

Appendix

For the sake of completeness and for future references, Tables 9-17 below contain detailed results of our computations for both (linear) UFL and qUFL.

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
MO1	100	100	1305.951410	1.83	1267.611394	1.47	2.94	45
MO2	100	100	1432.357320	1.66	1384.516719	1.42	3.34	32
MO3	100	100	1516.773000	2.76	1468.264089	2.44	3.20	47
MO4	100	100	1442.236430	1.44	1418.142416	1.29	1.67	10
MO5	100	100	1408.766380	1.59	1368.742892	1.29	2.84	33
MP1	200	200	2686.479460	5.40	2587.332224	4.41	3.69	58
MP2	200	200	2904.859010	5.25	2782.423128	3.37	4.21	76
MP3	200	200	2623.708880	2.74	2549.153223	2.33	2.84	26
MP4	200	200	2938.750020	9.29	2806.496012	7.54	4.50	142
MP5	200	200	2932.331110	7.74	2765.469658	5.47	5.69	272
MQ1	300	300	4091.009460	7.92	3943.844272	5.68	3.60	42
MQ2	300	300	4028.325730	11.20	3877.718770	7.50	3.74	119
MQ3	300	300	4275.431670	7.79	4128.671150	5.73	3.43	56
MQ4	300	300	4235.147010	7.76	4066.516145	5.35	3.98	72
MQ5	300	300	4080.742900	13.66	3860.071913	9.12	5.41	367
MR1	500	500	2608.148329	30.22	2438.570062	17.79	6.50	505
MR2	500	500	2654.734663	28.16	2517.943208	16.39	5.15	138
MR3	500	500	2788.250186	75.79	2592.583683	16.53	7.02	1069
MR4	500	500	2756.038599	49.88	2549.634566	15.38	7.49	1087
MR5	500	500	2505.047753	32.48	2344.396409	18.10	6.41	682
MS1	1000	1000	5283.757284	142.80	4930.109190	15.82	6.69	603
MT1	2000	2000	10069.802769	4889.77	9145.596419	41.11	9.18	10675

Table 6: All M instances for UFL solved to optimality (linear costs).

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
ga250a-1	250	250	257957	48.42	257621.391000	6.32	0.13	15408
ga250a-2	250	250	257502	16.21	257256.620721	7.27	0.10	3402
ga250a-3	250	250	257953	493.49	257554.773407	12.77	0.15	200184
ga250a-4	250	250	257987	94.90	257661.649685	8.96	0.13	31253
ga250a-5	250	250	258190	585.93	257790.245068	9.65	0.15	229446
ga250b-1	250	250	276296	459.35	273304.607325	9.57	1.08	79837
ga250b-2	250	250	275141	88.12	272725.591893	10.93	0.88	15382
ga250b-3	250	250	276093	294.31	273465.959876	13.01	0.95	44751
ga250b-4	250	250	276332	309.76	273674.972210	9.66	0.96	42305
ga250b-5	250	250	276404	239.05	273685.605552	11.31	0.98	34390
ga250c-1	250	250	334135	43.73	322972.561366	11.62	3.34	2022
ga250c-2	250	250	330728	29.89	321247.027668	12.86	2.87	766
ga250c-3	250	250	333662	37.86	322883.944855	13.17	3.23	1506
ga250c-4	250	250	332423	31.80	322223.157025	11.64	3.07	995
ga250c-5	250	250	333538	42.42	323060.428106	15.68	3.14	1406
gs250a-1	250	250	257964	44.41	257646.226693	7.75	0.12	15113
gs250a-2	250	250	257573	17.13	257295.404020	4.75	0.11	4160
gs250a-3	250	250	257626	143.30	257261.406140	6.24	0.14	56551
gs250a-4	250	250	257961	51.86	257612.040938	9.14	0.14	14474
gs250a-5	250	250	257896	105.20	257576.359942	14.76	0.12	35982
gs250b-1	250	250	276761	1553.22	273706.706855	11.62	1.10	274090
gs250b-2	250	250	275675	197.25	273039.140052	11.51	0.96	30610
gs250b-3	250	250	275710	265.83	273055.856146	11.46	0.96	37434
gs250b-4	250	250	276114	117.96	273744.809146	9.99	0.86	16871
gs250b-5	250	250	275916	170.12	273376.947727	10.71	0.92	25977
gs250c-1	250	250	332935	36.91	322714.264989	11.87	3.07	1064
gs250c-2	250	250	334630	55.20	323223.399647	17.60	3.41	2613
gs250c-3	250	250	333000	49.72	322004.934883	16.54	3.30	1623
gs250c-4	250	250	333158	34.35	322939.375070	13.18	3.07	1092
gs250c-5	250	250	334635	53.79	322663.655073	19.59	3.58	3279

Table 7: All g*250 instances for UFL solved to optimality (linear costs). Previously unknown optimal solutions shown in boldface.

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
ga500c-1	500	500	621360	4101.73	602860.484161	20.24	2.98	72942
ga500c-2	500	500	621464	8321.40	602811.779639	20.06	3.00	146412
ga500c-3	500	500	621428	12470.10	602379.469815	15.00	3.07	199101
ga500c-4	500	500	621754	9350.25	603259.772223	13.43	2.97	150378
ga500c-5	500	500	621313	9226.86	601500.282332	12.31	3.19	195191
gs500c-1	500	500	620041	7801.28	601985.824485	16.03	2.91	74511
gs500c-2	500	500	620434	4694.93	602299.966508	13.37	2.92	78969
gs500c-3	500	500	621204	11448.19	601980.526816	13.44	3.09	194657
gs500c-4	500	500	620437	7409.22	602136.181935	12.12	2.95	122315
gs500c-5	500	500	623180	26828.91	603115.401650	14.20	3.22	270147

Table 8: All g*500c instances for UFL solved to optimality (linear costs). Previously unknown optimal solutions shown in boldface.

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
B1.1	50	100	14032.005519	0.63	13696.521548	0.58	2.39	16
B1.2	50	100	13847.672000	0.36	13325.513467	0.25	3.77	49
B1.3	50	100	13388.124870	0.61	13214.726597	0.59	1.30	13
B1.4	50	100	12942.112566	0.92	12548.628511	0.83	3.04	32
B1.5	50	100	14207.725955	0.74	13860.584643	0.66	2.44	20
B1.6	50	100	12529.701053	0.73	12444.523201	0.71	0.68	0
B1.7	50	100	12781.344338	0.93	12550.229173	0.87	1.81	18
B1.8	50	100	12751.823705	0.67	12293.441936	0.56	3.59	43
B1.9	50	100	15674.462718	0.81	15203.088161	0.72	3.01	28
B1.10	50	100	13069.052954	0.56	12772.900090	0.49	2.27	24
C1.1	50	100	11000.439904	2.45	10368.787513	0.90	5.74	552
C1.2	50	100	10601.088985	1.82	9808.458064	0.74	7.48	539
C1.3	50	100	11330.445944	4.00	10657.286238	1.05	5.94	780
C1.4	50	100	10789.320016	2.00	10105.355062	0.85	6.34	394
C1.5	50	100	10967.008192	6.85	10241.066325	0.77	6.62	1895
C1.6	50	100	11296.490996	7.32	10514.043797	0.99	6.93	2330
C1.7	50	100	10980.997547	5.23	10292.306829	0.67	6.27	1654
C1.8	50	100	11058.606218	2.21	10469.183633	1.09	5.33	407
C1.9	50	100	11310.256519	4.13	10600.234180	0.63	6.28	1085
C1.10	50	100	11178.977525	3.28	10560.313565	0.80	5.53	737

Table 9: Performance of our slim qUFL model on B and C instances from ORLIB (quadratic costs).

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
D1.1	30	80	9284.373199	21.68	8531.717761	0.31	8.11	6723
D1.2	30	80	9237.865767	7.11	8559.684886	0.84	7.34	2836
D1.3	30	80	9185.572495	5.27	8627.387837	0.55	6.08	1203
D1.4	30	80	8783.343134	2.17	8155.361064	0.13	7.15	965
D1.5	30	80	8401.322106	0.38	7893.231472	0.10	6.05	153
D1.6	30	80	8988.476690	2.17	8480.254075	0.16	5.65	660
D1.7	30	80	9031.369277	2.86	8438.537164	0.18	6.56	977
D1.8	30	80	9000.747514	4.13	8325.166951	0.20	7.51	1606
D1.9	30	80	9048.394344	6.89	8401.782259	0.26	7.15	1897
D1.10	30	80	9208.919872	10.68	8582.566276	0.36	6.80	2862
D2.1	30	80	13021.521001	0.76	12004.378657	0.12	7.81	337
D2.2	30	80	13310.868193	3.72	12230.797282	0.16	8.11	1335
D2.3	30	80	13466.964045	3.14	12388.707324	0.19	8.01	1184
D2.4	30	80	13178.451378	2.92	12018.880338	0.18	8.80	1211
D2.5	30	80	13502.916063	4.06	12369.795737	0.16	8.39	1462
D2.6	30	80	13627.609516	3.90	12503.552184	0.18	8.25	1784
D2.7	30	80	13180.518073	0.84	12203.928546	0.13	7.41	324
D2.8	30	80	13484.419498	5.70	12277.824963	0.22	8.95	2324
D2.9	30	80	12964.500513	1.16	11788.404712	0.11	9.07	536
D2.10	30	80	12923.134799	1.68	11859.961446	0.13	8.23	831
D3.1	30	80	17026.170560	1.74	15396.699321	0.15	9.57	1137
D3.2	30	80	16848.371397	1.20	15424.453102	0.13	8.45	516
D3.3	30	80	16856.230221	0.97	15417.038262	0.10	8.54	464
D3.4	30	80	16388.369073	0.52	15081.064236	0.12	7.98	247
D3.5	30	80	16897.791661	1.45	15683.208172	0.17	7.19	403
D3.6	30	80	16536.742505	2.63	14908.235833	0.21	9.85	1069
D3.7	30	80	16603.090803	1.36	15301.280139	0.20	7.84	478
D3.8	30	80	16869.103468	1.56	15482.636643	0.17	8.22	572
D3.9	30	80	17158.999382	2.75	15614.544827	0.19	9.00	1229
D3.10	30	80	16394.007948	0.73	14987.718959	0.13	8.58	341
D4.1	30	80	20125.393247	2.29	18115.903476	0.15	9.98	1282
D4.2	30	80	19648.401006	0.85	17963.641987	0.10	8.57	424
D4.3	30	80	19904.577425	1.08	18395.187066	0.16	7.58	350
D4.4	30	80	20199.034055	1.35	18541.135732	0.13	8.21	611
D4.5	30	80	19349.624516	0.62	17870.799537	0.15	7.64	203
D4.6	30	80	19636.006145	0.73	17817.346871	0.10	9.26	424
D4.7	30	80	19350.866030	0.76	17714.257670	0.12	8.46	313
D4.8	30	80	20038.454210	0.98	18329.796590	0.16	8.53	556
D4.9	30	80	19554.801682	0.45	18035.086879	0.07	7.77	163
D4.10	30	80	19348.352767	0.42	17943.602055	0.09	7.26	156
D5.1	30	80	22888.887775	0.65	20946.639577	0.12	8.49	498
D5.2	30	80	23064.519598	0.69	21027.934505	0.10	8.83	434
D5.3	30	80	22627.595840	0.93	20545.200704	0.13	9.20	629
D5.4	30	80	22718.908473	0.70	20515.170664	0.10	9.70	560
D5.5	30	80	23022.942003	1.01	20778.991093	0.12	9.75	813
D5.6	30	80	21587.627368	0.38	19711.299984	0.10	8.69	239
D5.7	30	80	22989.638497	0.87	21173.076612	0.13	7.90	424
D5.8	30	80	22321.942170	0.49	20412.544074	0.11	8.55	393
D5.9	30	80	22919.303899	0.97	20743.544989	0.12	9.49	1048
D5.10	30	80	22461.918942	0.52	20671.296042	0.11	7.97	333

Table 10: Performance of our slim qUFL model on D instances from ORLIB (quadratic costs).

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	rootgap	nodes
D6.1	30	80	24678.920149	0.45	22505.961208	0.11	8.80	341
D6.2	30	80	25732.813136	1.03	23330.051121	0.18	9.34	772
D6.3	30	80	25122.264112	0.65	22579.264192	0.09	10.12	554
D6.4	30	80	25618.360201	1.17	23002.200197	0.14	10.21	1044
D6.5	30	80	26031.865385	1.00	23142.839464	0.13	11.10	1508
D6.6	30	80	25695.227321	1.03	23089.676676	0.13	10.14	1031
D6.7	30	80	25611.090028	0.78	22864.890972	0.09	10.72	796
D6.8	30	80	25442.084997	1.15	22863.173759	0.14	10.14	939
D6.9	30	80	25343.086094	1.01	22967.146724	0.14	9.38	610
D6.10	30	80	26223.032633	1.56	23670.539552	0.15	9.73	1281
D7.1	30	80	27002.443610	0.47	24675.731394	0.09	8.62	298
D7.2	30	80	27605.644559	0.56	25134.022297	0.11	8.95	403
D7.3	30	80	27386.712016	0.50	25230.915368	0.11	7.87	234
D7.4	30	80	28283.319555	0.86	25695.115422	0.13	9.15	778
D7.5	30	80	26403.466731	0.35	24731.858032	0.10	6.33	111
D7.6	30	80	27929.294395	0.77	25405.487965	0.13	9.04	620
D7.7	30	80	27326.000877	0.43	25153.435844	0.10	7.95	209
D7.8	30	80	26515.027216	0.22	24949.003833	0.08	5.91	71
D7.9	30	80	27589.745300	1.30	24765.889939	0.16	10.24	977
D7.10	30	80	27199.743532	0.55	24677.714521	0.14	9.27	389
D8.1	30	80	29141.169395	0.36	26951.454796	0.09	7.51	142
D8.2	30	80	29247.565012	0.41	27097.239568	0.10	7.35	124
D8.3	30	80	29825.561122	0.68	27224.999883	0.12	8.72	390
D8.4	30	80	29772.060194	0.33	27439.194733	0.08	7.84	191
D8.5	30	80	30660.545469	0.70	27971.526384	0.12	8.77	444
D8.6	30	80	29885.577120	0.59	27182.766531	0.11	9.04	397
D8.7	30	80	29679.568945	0.30	27064.552302	0.08	8.81	228
D8.8	30	80	29178.829180	0.32	26711.566384	0.09	8.46	155
D8.9	30	80	30192.659652	0.76	27172.532681	0.13	10.00	569
D8.10	30	80	29487.436510	0.27	27231.291497	0.09	7.65	191
D9.1	30	80	32410.989472	0.40	29962.436486	0.09	7.55	211
D9.2	30	80	30624.047673	0.24	28404.421095	0.08	7.25	86
D9.3	30	80	31347.453645	0.43	29043.508623	0.10	7.35	115
D9.4	30	80	31982.973804	0.40	29516.771505	0.10	7.71	192
D9.5	30	80	31417.137240	0.26	28975.939885	0.06	7.77	129
D9.6	30	80	31667.617861	0.27	29265.235771	0.09	7.59	144
D9.7	30	80	31198.246667	0.68	28813.150827	0.15	7.64	188
D9.8	30	80	32391.068672	0.68	29531.715443	0.15	8.83	524
D9.9	30	80	31193.617388	0.27	29196.222413	0.10	6.40	117
D9.10	30	80	31697.508219	0.48	29269.083360	0.10	7.66	267
D10.1	30	80	33554.380582	0.39	31128.189530	0.12	7.23	185
D10.2	30	80	33350.621280	0.33	30822.613884	0.10	7.58	188
D10.3	30	80	31872.225127	0.24	30391.708179	0.10	4.65	27
D10.4	30	80	33774.768267	0.34	31219.318070	0.11	7.57	205
D10.5	30	80	33040.626061	0.32	30760.626431	0.10	6.90	118
D10.6	30	80	33411.062732	0.38	31130.271189	0.09	6.83	115
D10.7	30	80	33347.673651	0.46	31075.721620	0.14	6.81	133
D10.8	30	80	34067.144655	0.52	31420.806099	0.14	7.77	244
D10.9	30	80	34386.529487	0.70	31619.967510	0.15	8.05	396
D10.10	30	80	33227.195483	0.30	30988.861190	0.08	6.74	101

Table 11: Performance of our slim qUFL model on D instances from ORLIB (quadratic costs).

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
E1.1	50	100	9929.156583	209.09	9156.177668	0.52	7.78	664692
E1.2	50	100	9742.682326	80.89	8903.187316	0.62	8.62	21894
E1.3	50	100	9866.224109	77.25	9069.251206	0.50	8.08	9781
E1.4	50	100	9444.198122	12.54	8658.017009	0.43	8.32	2483
E1.5	50	100	9537.061284	40.12	8771.527942	0.55	8.03	4976
E1.6	50	100	9503.301403	14.48	8896.167086	0.49	6.39	1526
E1.7	50	100	9630.576179	27.36	8870.376248	0.49	7.89	4496
E1.8	50	100	9422.206703	15.43	8682.932771	0.45	7.85	2995
E1.9	50	100	9644.957935	13.99	8952.729258	0.41	7.18	1881
E1.10	50	100	9670.795716	442.15	8746.101228	0.47	9.56	1853869
E2.1	50	100	14534.587758	24.72	13433.681971	0.43	7.57	3020
E2.2	50	100	14335.112341	19.67	13188.032580	0.50	8.00	2135
E2.3	50	100	14486.774849	25.05	13129.714027	0.40	9.37	3844
E2.4	50	100	14553.684124	9.37	13404.539309	0.40	7.90	1354
E2.5	50	100	14267.719588	31.58	13030.586416	0.55	8.67	3429
E2.6	50	100	14611.020505	55.43	13362.722269	0.43	8.54	6587
E2.7	50	100	14624.495393	19.58	13405.151738	0.36	8.34	2730
E2.8	50	100	14587.225829	28.89	13262.099961	0.40	9.08	4490
E2.9	50	100	14556.708968	30.53	13215.493845	0.42	9.21	5002
E2.10	50	100	14860.081534	44.09	13428.857525	0.52	9.63	8206
E3.1	50	100	18313.542127	12.09	16539.942534	0.49	9.68	2817
E3.2	50	100	18539.324556	16.36	16793.331388	0.33	9.42	3237
E3.3	50	100	18519.084064	23.94	16763.985185	0.35	9.48	3855
E3.4	50	100	18587.284615	25.29	16883.920170	0.40	9.16	4459
E3.5	50	100	18536.855073	50.40	16642.114865	0.43	10.22	8464
E3.6	50	100	18408.644369	54.45	16378.232605	0.43	11.03	13240
E3.7	50	100	18504.032907	17.85	16626.283525	0.42	10.15	4800
E3.8	50	100	18192.988827	9.81	16546.078963	0.36	9.05	2215
E3.9	50	100	18257.866537	19.93	16381.046460	0.34	10.28	4028
E3.10	50	100	18255.735669	8.21	16409.589478	0.31	10.11	2805
E4.1	50	100	21938.856834	16.31	19914.977413	0.44	9.23	2800
E4.2	50	100	20577.661850	2.24	18887.660669	0.26	8.21	360
E4.3	50	100	20891.425609	1.79	19169.206800	0.32	8.24	410
E4.4	50	100	21257.474682	5.00	19425.219820	0.33	8.62	851
E4.5	50	100	21458.851100	9.17	19475.430435	0.33	9.24	1659
E4.6	50	100	21728.120485	6.18	19902.316345	0.30	8.40	1116
E4.7	50	100	21793.987030	13.63	19560.738579	0.41	10.25	3235
E4.8	50	100	21507.933737	8.45	19426.331399	0.36	9.68	2025
E4.9	50	100	21625.919329	12.41	19452.080523	0.35	10.05	2314
E4.10	50	100	21136.577294	4.97	19141.688805	0.30	9.44	1077
E5.1	50	100	24318.745489	6.80	21994.897970	0.38	9.56	1634
E5.2	50	100	24149.635771	2.98	21874.182446	0.26	9.42	851
E5.3	50	100	24895.332124	8.27	22589.645681	0.37	9.26	1944
E5.4	50	100	24327.653048	2.38	22196.098205	0.26	8.76	676
E5.5	50	100	24797.810470	7.28	22387.803887	0.29	9.72	2029
E5.6	50	100	24713.466670	6.57	22434.911632	0.33	9.22	1683
E5.7	50	100	23955.752674	1.53	21815.393350	0.21	8.93	369
E5.8	50	100	24826.221438	5.46	22624.773688	0.32	8.87	1055
E5.9	50	100	24662.629818	4.66	22554.089607	0.33	8.55	806
E5.10	50	100	24667.232023	5.32	22267.057816	0.33	9.73	1423

Table 12: Performance of our slim qUFL model on E instances from ORLIB (quadratic costs).

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
E6.1	50	100	27245.865153	3.21	24574.546091	0.33	9.80	1089
E6.2	50	100	27618.608596	5.87	25131.504571	0.32	9.01	1694
E6.3	50	100	28093.999469	7.03	25305.643966	0.26	9.93	2326
E6.4	50	100	27831.991802	8.99	24870.412030	0.34	10.64	2919
E6.5	50	100	27105.307497	3.30	24499.136430	0.27	9.61	1240
E6.6	50	100	27680.075708	13.38	24889.993643	0.45	10.08	3010
E6.7	50	100	28365.709173	12.57	25513.847453	0.43	10.05	4619
E6.8	50	100	27305.100107	3.34	24661.508366	0.26	9.68	1149
E6.9	50	100	27577.795348	6.91	24628.499479	0.31	10.69	2139
E6.10	50	100	27980.777074	5.47	25147.654422	0.27	10.13	2076
E7.1	50	100	30862.777552	6.44	27483.703051	0.33	10.95	4803
E7.2	50	100	30728.535541	10.85	27250.593590	0.39	11.32	5135
E7.3	50	100	30521.324742	6.45	26752.800427	0.31	12.35	5841
E7.4	50	100	30680.614284	5.65	27392.980867	0.35	10.72	3206
E7.5	50	100	31240.979424	12.15	27775.354648	0.33	11.09	6077
E7.6	50	100	30504.794201	10.80	26770.014191	0.40	12.24	5698
E7.7	50	100	30487.241879	4.93	27217.882581	0.29	10.72	2220
E7.8	50	100	29699.466786	2.08	27091.657461	0.29	8.78	661
E7.9	50	100	30377.444112	5.41	27149.156563	0.25	10.63	2428
E7.10	50	100	30962.017491	8.89	27451.003328	0.30	11.34	4614
E8.1	50	100	32441.820572	3.11	29053.431034	0.29	10.44	1217
E8.2	50	100	32796.427832	3.46	29563.379136	0.29	9.86	1350
E8.3	50	100	32610.486601	3.05	29131.722520	0.29	10.67	1852
E8.4	50	100	32481.548094	2.17	29570.629612	0.28	8.96	691
E8.5	50	100	32870.424399	4.11	29222.882972	0.32	11.10	2007
E8.6	50	100	32632.269750	2.93	29241.030072	0.31	10.39	1225
E8.7	50	100	33620.975452	9.59	29640.958711	0.32	11.84	7099
E8.8	50	100	32691.540818	3.97	29121.385089	0.28	10.92	1733
E8.9	50	100	34038.569387	8.13	30063.528729	0.31	11.68	9186
E8.10	50	100	32763.642530	2.78	29306.571503	0.27	10.55	1878
E9.1	50	100	35385.680159	2.92	31508.104162	0.26	10.96	2331
E9.2	50	100	35246.889116	3.73	32080.466853	0.38	8.98	1170
E9.3	50	100	34733.180932	2.07	31522.101142	0.29	9.24	790
E9.4	50	100	34593.013200	2.27	31589.879797	0.40	8.68	777
E9.5	50	100	34911.205798	2.08	31744.425987	0.35	9.07	1022
E9.6	50	100	33992.013923	0.90	30720.323308	0.18	9.62	521
E9.7	50	100	34737.834473	1.93	31213.104096	0.30	10.15	1195
E9.8	50	100	35450.469761	3.29	32073.927413	0.33	9.52	1369
E9.9	50	100	35073.964306	3.32	31601.538240	0.31	9.90	1512
E9.10	50	100	34632.709909	2.26	31635.521968	0.31	8.65	715
E10.1	50	100	36844.713067	2.27	33655.589959	0.32	8.66	646
E10.2	50	100	37526.691156	2.27	33927.429038	0.29	9.59	1264
E10.3	50	100	36689.234782	1.73	33344.844749	0.34	9.12	532
E10.4	50	100	36627.580809	1.67	33419.329042	0.27	8.76	486
E10.5	50	100	37838.851838	3.38	34486.239254	0.27	8.86	1004
E10.6	50	100	36330.698056	0.88	33595.452997	0.23	7.53	262
E10.7	50	100	37075.482787	2.24	33769.623091	0.31	8.92	742
E10.8	50	100	36686.460103	1.43	33554.653678	0.22	8.54	349
E10.9	50	100	37382.126357	2.14	34210.942526	0.32	8.48	736
E10.10	50	100	37813.965988	3.42	34106.139690	0.40	9.81	1579

Table 13: Performance of our slim qUFL model on E instances from ORLIB (quadratic costs).

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
cap101.txt	25	50	198198.872443	0.01	198198.872444	0.01	-0.00	0
cap102.txt	25	50	249809.185072	0.01	249808.125659	0.01	0.00	0
cap103.txt	25	50	290945.894297	0.03	290571.850341	0.01	0.13	4
cap104.txt	25	50	341163.118782	0.02	341121.871252	0.02	0.01	0
cap111.txt	50	50	196692.072237	0.02	196682.095733	0.02	0.01	3
cap112.txt	50	50	248407.264312	0.03	248272.090741	0.02	0.05	5
cap113.txt	50	50	289641.492435	0.04	289423.224848	0.02	0.08	6
cap114.txt	50	50	340489.888416	0.03	340234.698044	0.02	0.07	5
cap121.txt	50	50	196692.072237	0.02	196682.095733	0.02	0.01	3
cap122.txt	50	50	248407.264312	0.03	248272.090741	0.02	0.05	5
cap123.txt	50	50	289641.492435	0.04	289423.224848	0.02	0.08	6
cap124.txt	50	50	340489.888416	0.03	340234.698044	0.02	0.07	5
cap131.txt	50	50	196692.072237	0.02	196682.095733	0.02	0.01	3
cap132.txt	50	50	248407.264312	0.03	248272.090741	0.02	0.05	5
cap133.txt	50	50	289641.492435	0.04	289423.224848	0.02	0.08	6
cap134.txt	50	50	340489.888416	0.03	340234.698044	0.02	0.07	5
cap41.txt	16	50	209147.513131	0.02	209078.381192	0.01	0.03	6
cap42.txt	16	50	258384.012313	0.01	258218.353676	0.00	0.06	0
cap43.txt	16	50	297469.742386	0.01	297440.144593	0.00	0.01	0
cap44.txt	16	50	346437.300838	0.01	346432.943779	0.00	0.00	0
cap51.txt	16	50	297469.742386	0.01	297440.144593	0.00	0.01	0
cap61.txt	16	50	209147.513131	0.02	209078.381192	0.01	0.03	6
cap62.txt	16	50	258384.012313	0.01	258218.353676	0.00	0.06	0
cap63.txt	16	50	297469.742386	0.01	297440.144593	0.00	0.01	0
cap64.txt	16	50	346437.300838	0.01	346432.943779	0.00	0.00	0
cap71.txt	16	50	209147.513131	0.02	209078.381192	0.01	0.03	6
cap72.txt	16	50	258384.012313	0.01	258218.353676	0.00	0.06	0
cap73.txt	16	50	297469.742386	0.01	297440.144593	0.00	0.01	0
cap74.txt	16	50	346437.300838	0.01	346432.943779	0.00	0.00	0
cap81.txt	25	50	198198.872443	0.01	198198.872444	0.01	-0.00	0
cap82.txt	25	50	249809.185072	0.01	249808.125659	0.01	0.00	0
cap83.txt	25	50	290945.894297	0.03	290571.850341	0.01	0.13	4
cap84.txt	25	50	341163.118782	0.02	341121.871252	0.02	0.01	0
cap91.txt	25	50	198198.872443	0.01	198198.872444	0.01	-0.00	0
cap92.txt	25	50	249809.185072	0.01	249808.125659	0.01	0.00	0
cap93.txt	25	50	290945.894297	0.03	290571.850341	0.01	0.13	4
cap94.txt	25	50	341163.118782	0.02	341121.871252	0.02	0.01	0
capa.txt	100	1000	11060090.523934	2.59	10898396.110654	1.62	1.46	71
capb.txt	100	1000	6891132.403879	4.24	6803190.552973	2.08	1.28	286
capc.txt	100	1000	5777026.274077	3.69	5726328.020829	2.27	0.88	99

Table 14: Performance of our slim qUFL model on cap* instances from ORLIB (quadratic costs).

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
gs250a-1	250	250	12633.858555	0.21	12633.271955	0.19	0.00	7
gs250a-2	250	250	12723.318779	0.20	12723.218469	0.19	0.00	0
gs250a-3	250	250	12636.391724	0.24	12636.338503	0.22	0.00	5
gs250a-4	250	250	12628.620604	0.20	12628.439925	0.19	0.00	5
gs250a-5	250	250	12604.942212	0.21	12604.783766	0.20	0.00	0
gs250b-1	250	250	38625.725625	0.24	38618.122355	0.21	0.02	21
gs250b-2	250	250	38934.394122	0.25	38928.793683	0.21	0.01	18
gs250b-3	250	250	38675.929807	0.21	38672.148776	0.19	0.01	0
gs250b-4	250	250	38777.891611	0.41	38761.526533	0.24	0.04	339
gs250b-5	250	250	38719.501675	0.25	38708.792989	0.22	0.03	25
gs250c-1	250	250	120474.797667	0.52	120439.921239	0.48	0.03	13
gs250c-2	250	250	122031.956567	0.49	121981.751982	0.44	0.04	30
gs250c-3	250	250	120863.499762	0.45	120819.787084	0.42	0.04	13
gs250c-4	250	250	121217.929629	0.34	121188.172224	0.30	0.02	11
gs250c-5	250	250	121316.362556	0.53	121273.396543	0.46	0.04	41
gs500a-1	500	500	17578.186397	0.82	17578.122078	0.77	0.00	10
gs500a-2	500	500	17645.410713	0.81	17645.316518	0.77	0.00	10
gs500a-3	500	500	17510.760092	0.85	17510.636113	0.76	0.00	35
gs500a-4	500	500	17529.727138	0.78	17529.663411	0.77	0.00	1
gs500a-5	500	500	17654.985621	0.80	17654.943824	0.76	0.00	6
gs500b-1	500	500	54221.008053	2.48	54218.481084	2.46	0.00	0
gs500b-2	500	500	54641.496690	2.68	54626.975702	2.61	0.03	21
gs500b-3	500	500	54385.233438	1.85	54349.634874	1.55	0.07	126
gs500b-4	500	500	54290.557516	3.53	54285.146373	3.23	0.01	190
gs500b-5	500	500	54539.325360	1.82	54509.909369	1.69	0.05	27
gs500c-1	500	500	170495.752651	21.97	170239.543935	1.23	0.15	14156
gs500c-2	500	500	171044.987106	2.45	170860.982585	1.18	0.11	656
gs500c-3	500	500	171236.065342	27.28	170987.933058	1.54	0.14	35441
gs500c-4	500	500	170740.764114	16.58	170492.376791	0.97	0.15	9880
gs500c-5	500	500	171170.211422	6.95	170905.040422	1.36	0.15	3475
gs750a-1	750	750	21481.746742	1.97	21480.716490	1.71	0.00	35
gs750a-2	750	750	21395.527255	2.20	21394.215388	1.64	0.01	143
gs750a-3	750	750	21484.689396	1.90	21483.591903	1.64	0.01	34
gs750a-4	750	750	21426.082803	1.94	21424.852668	1.64	0.01	51
gs750a-5	750	750	21488.509512	1.68	21488.358552	1.64	0.00	3
gs750b-1	750	750	66621.039307	1.97	66619.107840	1.82	0.00	29
gs750b-2	750	750	66463.658421	8.54	66456.042264	1.82	0.01	2011
gs750b-3	750	750	66734.171061	1.82	66733.713927	1.79	0.00	0
gs750b-4	750	750	66626.128388	2.00	66623.377493	1.84	0.00	28
gs750b-5	750	750	66506.140140	1.86	66501.966024	1.83	0.01	2
gs750c-1	750	750	209089.008287	19.68	208946.188003	2.39	0.07	5933
gs750c-2	750	750	209023.548416	36.53	208870.169230	3.03	0.07	19103
gs750c-3	750	750	209357.291160	10.46	209215.563124	4.47	0.07	2109
gs750c-4	750	750	209438.506671	30.04	209286.575163	2.80	0.07	7446
gs750c-5	750	750	209229.646295	37.98	209073.658210	2.23	0.07	49594

Table 15: Performance of our slim qUFL model on gs* instances from UFLLIB (quadratic costs).

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
ga250a-1	250	250	12719.533977	0.428786	12719.530507	0.415388	0.000027	0
ga250a-2	250	250	12552.936248	0.456415	12552.849390	0.443985	0.000692	0
ga250a-3	250	250	12701.173983	0.465590	12700.554603	0.445361	0.004877	9
ga250a-4	250	250	12596.578129	0.491309	12596.111277	0.465797	0.003706	13
ga250a-5	250	250	12619.414974	0.173889	12619.313419	0.162894	0.000805	0
ga250b-1	250	250	38836.625923	0.220798	38832.838550	0.206783	0.009752	0
ga250b-2	250	250	38457.705301	0.278533	38441.757806	0.230750	0.041468	38
ga250b-3	250	250	38854.416851	0.245665	38849.804970	0.208089	0.011870	25
ga250b-4	250	250	38818.594279	0.267422	38805.975132	0.216957	0.032508	58
ga250b-5	250	250	38660.119334	0.371695	38643.367602	0.240501	0.043331	235
ga250c-1	250	250	121804.946988	0.329247	121794.550774	0.318951	0.008535	3
ga250c-2	250	250	120202.866912	0.297556	120158.102300	0.267686	0.037241	18
ga250c-3	250	250	121328.316329	0.340634	121283.752051	0.296418	0.036730	38
ga250c-4	250	250	121831.490417	0.524049	121778.886045	0.442793	0.043178	93
ga250c-5	250	250	121081.938730	0.517018	121023.722668	0.473089	0.048080	45
ga500a-1	500	500	17728.229656	0.759997	17728.105721	0.736688	0.000699	0
ga500a-2	500	500	17601.647513	0.759103	17601.555149	0.736135	0.000525	0
ga500a-3	500	500	17628.411284	0.801732	17628.151470	0.762048	0.001474	7
ga500a-4	500	500	17608.899247	0.771241	17608.726077	0.731291	0.000983	8
ga500a-5	500	500	17553.849070	0.731572	17553.768482	0.704281	0.000459	0
ga500b-1	500	500	54617.814966	1.773924	54592.516651	1.714924	0.046319	2
ga500b-2	500	500	54309.398377	1.919447	54280.401192	1.804601	0.053393	13
ga500b-3	500	500	54510.983133	3.179420	54506.301100	3.056560	0.008589	61
ga500b-4	500	500	54453.425408	2.008263	54421.595522	1.584899	0.058453	196
ga500b-5	500	500	54248.375667	1.731638	54223.277515	1.604858	0.046265	18
ga500c-1	500	500	171440.595530	9.849523	171189.753151	1.443641	0.146314	3998
ga500c-2	500	500	170966.289717	31.460599	170648.578446	1.992068	0.185833	135122
ga500c-3	500	500	171064.595331	4.383382	170827.714508	1.031294	0.138474	2325
ga500c-4	500	500	171355.977333	25.158046	171051.316116	1.997193	0.177794	87276
ga500c-5	500	500	170732.019561	26.468147	170445.056411	1.002103	0.168078	20837
ga750a-1	750	750	21454.097467	1.796982	21453.272117	1.628245	0.003847	15
ga750a-2	750	750	21312.186688	2.342660	21310.706487	1.634806	0.006945	195
ga750a-3	750	750	21343.820134	2.702123	21342.737437	1.625286	0.005073	324
ga750a-4	750	750	21540.541353	1.685356	21540.473492	1.629486	0.000315	3
ga750a-5	750	750	21294.696772	1.629846	21294.669133	1.591675	0.000130	0
ga750b-1	750	750	66457.646067	1.978090	66455.331052	1.834522	0.003483	26
ga750b-2	750	750	66189.378762	2.321838	66183.596058	1.831737	0.008737	144
ga750b-3	750	750	66295.153760	2.175801	66290.530620	1.830254	0.006974	79
ga750b-4	750	750	66658.392142	1.808612	66656.073961	1.760198	0.003478	3
ga750b-5	750	750	66409.403295	2.131812	66403.465973	1.824355	0.008940	74
ga750c-1	750	750	208746.415344	36.124844	208560.918557	2.597722	0.088862	26329
ga750c-2	750	750	208328.530052	39.892769	208171.648820	2.236358	0.075305	28426
ga750c-3	750	750	208625.496768	36.297678	208489.573312	2.957719	0.065152	13775
ga750c-4	750	750	209338.346423	18.939491	209162.272345	2.368603	0.084110	5100
ga750c-5	750	750	209122.698081	47.718554	208947.757154	1.868571	0.083655	248060

Table 16: Performance of our slim qUFL model on ga* instances from UFLLIB (quadratic costs).

inst.	n	m	opt	$t[s]$	rootbound	$t_{root}[s]$	$g_r[\%]$	nodes
MO1	100	100	618.042992	0.09	617.279304	0.08	0.12	4
MO2	100	100	667.118745	0.09	666.543774	0.08	0.09	3
MO3	100	100	716.883255	0.09	710.594541	0.06	0.88	12
MO4	100	100	680.847385	0.10	679.902555	0.09	0.14	3
MO5	100	100	714.424108	0.09	714.020554	0.08	0.06	2
MP1	200	200	1279.256448	0.29	1275.587063	0.25	0.29	10
MP2	200	200	1416.571645	0.35	1402.774815	0.28	0.97	22
MP3	200	200	1256.660588	0.30	1254.901421	0.28	0.14	1
MP4	200	200	1376.371490	0.24	1373.406966	0.21	0.22	5
MP5	200	200	1398.789178	0.28	1394.092315	0.26	0.34	5
MQ1	300	300	1950.890622	0.47	1947.433944	0.46	0.18	1
MQ2	300	300	1851.116781	0.55	1846.760328	0.49	0.24	8
MQ3	300	300	2061.005740	0.61	2046.729810	0.51	0.69	16
MQ4	300	300	2067.886319	0.54	2059.745723	0.47	0.39	9
MQ5	300	300	1892.504952	0.65	1888.245346	0.58	0.23	10
MR1	500	500	1289.878712	3.09	1277.654170	2.60	0.95	36
MR2	500	500	1420.464199	3.48	1401.644119	2.82	1.32	58
MR3	500	500	1431.199847	3.56	1415.799274	2.90	1.08	39
MR4	500	500	1402.027193	2.43	1393.565153	2.13	0.60	16
MR5	500	500	1254.645421	2.36	1246.912227	2.07	0.62	17
MS1	1000	1000	2672.562217	8.28	2664.219242	7.86	0.31	6
MT1	2000	2000	5133.779083	81.99	5075.996465	65.24	1.13	165

Table 17: Performance of our slim qUFL model on M instances from UFLLIB (quadratic costs).