



CIRRELT

Centre interuniversitaire de recherche
sur les réseaux d'entreprise, la logistique et le transport

Interuniversity Research Centre
on Enterprise Networks, Logistics and Transportation

A Benders Decomposition Approach for the Symmetric TSP with Generalized Latency

Fausto Errico
Teodor Gabriel Crainic
Federico Malucelli
Maddalena Nonato

December 2012

CIRRELT-2012-78

Bureaux de Montréal :

Université de Montréal
C.P. 6128, succ. Centre-ville
Montréal (Québec)
Canada H3C 3J7
Téléphone : 514 343-7575
Télécopie : 514 343-7121

Bureaux de Québec :

Université Laval
2325, de la Terrasse, bureau 2642
Québec (Québec)
Canada G1V 0A6
Téléphone : 418 656-2073
Télécopie : 418 656-2624

www.cirrelt.ca

A Benders Decomposition Approach for the Symmetric TSP with Generalized Latency

Fausto Errico^{1,2}, Teodor Gabriel Crainic^{1,3,*}, Federico Malucelli⁴, Maddalena Nonato⁵

¹ Interuniversity Research Centre on Enterprise Networks, Logistics and Transportation (CIRRELT)

² Department of Mathematics and Industrial Engineering, École Polytechnique de Montréal, P.O. Box 6079, Station Centre-ville, Montréal, Canada H3C 3A7

³ Department of Management and Technology, Université du Québec à Montréal, P.O. Box 8888, Station Centre-Ville, Montréal, Canada H3C 3P8

⁴ Dipartimento di Elettronica ed Informazione, Politecnico di Milano, Piazza L. Da Vinci 32, 20133 Milano, Italy

⁵ Dipartimento di Ingegneria, Università degli studi di Ferrara, Via Saragat 1, 44100, Ferrara, Italy

Abstract. We present a new problem within the class of TSP. The problem arises in the transportation field, namely in the design of semi-flexible transit lines, but it may have applications also in other areas. The objective function combines the usual cost of the Hamiltonian circuit with a routing cost that depends on the actual flow on each arc. The main contributions include the formulation in terms of multicommodity flow, the study of mathematical properties and a Branch-and-Cut approach based on Benders decomposition. Some variants of the algorithm are experimentally evaluated, in comparison with a commercial solver. Several algorithmic refinements are proposed and their efficacy is experimentally tested on a variety of instances. Tests on benchmark instances confirm the validity of our approach.

Keywords: Traveling salesperson problem, latency, Benders decomposition, branch-and-cut, semi-flexible transit.

Acknowledgements. While working on this paper, F. Errico was postdoctoral researcher at École des sciences de la gestion, Université du Québec à Montréal. T.G. Crainic was the NSERC Industrial Research Chair in Logistics Management, ESG UQAM, and Adjunct Professor with the Department of Computer Science and Operations Research, Université de Montréal, and the Department of Economics and Business Administration, Molde University College, Norway, and F. Errico was postdoctoral researcher with the Chair. Funding for this project was provided by the Natural Sciences and Engineering Research Council of Canada (NSERC), through its Industrial Research Chair and Discovery Grant programs, by our partners CN, Rona, Alimentation Couche-Tard, the ministère des Transports du Québec, and by the Fonds de recherche du Québec - Nature et technologies (FRQNT) through its Team Research Project program. We also gratefully acknowledge the support of the Canada Foundation for Innovation, the Québec Ministry of Education, UQAM and SUN/Azuris for our computing infrastructure, as well as that of Fonds de recherche du Québec for the CIRRELT infrastructure grants.

Results and views expressed in this publication are the sole responsibility of the authors and do not necessarily reflect those of CIRRELT.

Les résultats et opinions contenus dans cette publication ne reflètent pas nécessairement la position du CIRRELT et n'engagent pas sa responsabilité.

* Corresponding author: TeodorGabriel.Crainic@cirrelt.ca

1 Introduction

In this paper, we introduce a new optimization problem, the TSP-GL, based on the Symmetric Traveling Salesperson Problem (S-TSP) to which we add a generalized latency component in the objective function. This problem has immediate application in the network design of transit lines but it may also be applied in other fields.

The problem can be described as follows. As in the S-TSP, we are given a complete undirected graph with non negative costs on the edges and a Hamiltonian circuit has to be found. In addition, we are given an origin/destination demand matrix, where each entry specifies the amount of demand to be routed between the corresponding origin and destination nodes of the graph. Moreover, we are given the traversing times or unit routing costs for every edge and for both directions. Therefore, besides finding a Hamiltonian circuit, we must also find the routing of the demand on the circuit. The aspect that characterizes the TSP-GL and makes the problem very challenging is the objective function that accounts for two components. The first one considers the cost of the edges in the Hamiltonian circuit (*installation cost*), the second accounts for the overall time spent to traverse the required portion of the circuit by each demand unit (*routing cost* or *latency*).

In the design of public transit services, passenger travel times constitute a key factor directly affecting the level of service. Traditionally, passenger travel times are accounted for at tactical planning level, i.e., after the line design has been established. Passenger travel times optimization is only considered when defining service frequencies and timetables (see Desaulniers and Hickman 2007, for an extensive review on public transit planning). In such a context, the TSP-GL plays an important role, allowing for the simultaneous optimization of passenger travel times (latency) and operating costs (installation).

For illustration sake, let us briefly describe a basic version of the design of a single circular transit line, operated in both directions, in terms of TSP-GL optimizing both operational costs and passenger travel times. The nodes of the graph correspond to the bus stops, the edges are the possible direct connections between two stops. The installation costs and the routing costs are given by the traveling cost and shortest traveling time between the corresponding pair of bus stops, respectively. The origin/destination matrix is the estimated transportation demand between bus stop pairs. In planning the line we must design a circuit passing by all the bus stops exactly once. The objective is to minimize a function that combines the traveling cost (operational cost) and the overall time spent by users on vehicles (level of service). Given a circuit, for any origin there are two possible alternative itineraries to reach the destination, depending on the direction taken on the circuit. We assume that users choose the shortest one in terms of traveling time. The choice of considering a circular line does not limit the generality of the approach, that can be easily adapted to Hamiltonian paths.

The TSP-GL has been applied to the design of a wide class of on-demand transit services called *semi-flexible*. Errico et al. (2011b), presents an extensive computational study of the application of TSP-GL to the design of semi-flexible transit lines, evaluating the influence of the two components on the final results. For an exhaustive review on the state-of-the-art and state-of-the-practice of semi-flexible transit systems see Errico et al. (2011a) and Potts et al. (2010).

The TSP-GL is at least as difficult as the TSP. In fact, if the latency component is removed from the objective function, the TSP-GL reduces to a S-TSP, which immediately proves the complexity of the problem. However, the TSP-GL is much harder than the TSP due to the particular objective function. Note that, given a Hamiltonian circuit, determining the routing of each demand pair is straightforward since we must only select in which direction it has to be sent along the circuit, and this can be done separately and independently for each demand pair. In spite of this fact, the nature of the objective function makes the

problem decisions extremely correlated and computationally very challenging, as we will see further on.

At the best of our knowledge, we are the first to introduce and formalize the TSP-GL. However, there are two families of problems known in the literature which are related to the TSP-GL: The Fixed-Charge Network Design Problem (FNDP) and several variants of the TSP. Nonetheless, the specific structure of the TSP-GL and its complexity do not allow for a successful adaptation of existing approaches to the above mentioned problems and actually, new ad hoc solution methods are required in order to solve practical size instances.

We propose a multicommodity flow based formulation for the TSP-GL as an extension of the traditional S-TSP formulation where the new variables are used to express our specific objective function. Such a formulation has the advantage of being compact, having a polynomial number of variables and constraints with respect to the dimension of the problem. However, this formulation presents some drawbacks: In spite of its compactness, the solution of the LP relaxation requires very long computing times, and the lower bounds obtained this way are quite loose. These reasons gave us the motivation to investigate a solution approach that combines Benders decomposition and generation of valid inequalities. The Benders decomposition, in fact, requires to iteratively solve a *master problem*, which is a relaxation of the original problem, and a (dual) *subproblem* used to generate inequalities improving the current solution of the master problem. The inequalities generated are of two kinds: *optimality* cuts and *feasibility* cuts.

Differently from the traditional scheme where Benders cuts are only generated at integer master solutions, in this paper we consider an alternative scheme where integrality constraints of the master problem are relaxed, and integrality is recovered by embedding the Benders master problem into a genuine Branch and Cut framework. We call such an approach the Benders-Branch-and-Cut algorithm (BBC). The BBC algorithm seems particularly attractive for the TSP-GL because, by projecting flow variables and relaxing integrality, the similarity between TSP-GL and S-TSP can be effectively exploited.

Starting from the traditional Benders reformulation, we consider three alternative families of inequalities to be deployed as feasibility cuts. We prove mathematical relations among the lower bounds provided by each family. We implemented the corresponding versions of the BBC algorithm and experimentally compared their efficiency. All BBC algorithms widely outperform a state-of-the-art MIP solver applied to the initial multicommodity flow formulation. As a second step, we refined the most efficient BBC algorithm by implementing some techniques known to improve the performance of the traditional Benders decomposition. We experimentally compare such refinements and the resulting best choices are used to attack a set of newly created benchmark instances derived from the well-known TSPLIB library. It turns out that the refined BBC algorithm is able to optimally solve instances with up to 42 nodes, while in the most difficult 101-node instance, it provides a feasible solution within 8.8% of the optimality in 5 hours of computing time.

The rest of the paper is organized as follows: In Section 2, we introduce the TSP-GL multicommodity flow based formulation and we survey the related works. In Section 3, we recall the main ideas of the Benders decomposition and we apply Benders reformulation to the TSP-GL. Section 4 describes the general BBC and motivates our choice. In the same section, three versions of the BBC algorithm and their mathematical properties are examined. In Section 5, we report the details of the experimental campaign, computational results and analyses. We draw our conclusions in Section 6.

2 Problem statement and related works

The TSP-GL problem involves two levels of decisions. On the one hand the set of edges defining the Hamiltonian circuit must be selected, on the other hand the demand has to be routed on the circuit and the routing direction is an issue. This suggests the employment of a mixed graph since both edges and directed

arcs are used.

Consider a complete mixed graph $G = (N, E \cup A)$, where $N = \{1, \dots, n\}$ is the set of nodes, $E = \{[i, j] : i, j \in N, i < j\}$ is the set of edges, and $A = \{(i, j) : i, j \in N, i \neq j\}$ is the set of directed arcs. A *design* cost $\bar{c}_e > 0$ is associated to each edge $e = [i, j] \in E$. Moreover, the amount of demand d_{hk} is given for each origin and destination node pair (h, k) . Let D denote the set of pairs having non zero demand, that is $D = \{(h, k) : h, k \in N, d_{hk} > 0\}$. A routing cost $q_{ij}^{hk} \geq 0$ is associated to every arc $(i, j) \in A$ and every demand $(h, k) \in D$.

The problem consists in finding a subset $X \subset E$ of edges defining a Hamiltonian circuit in G and a flow such that the demand is satisfied, that is the flow going from h to k is equal to d_{hk} for each pair of nodes $(h, k) \in D$ and the flow is consistent with the circuit definition, that is the flow on arc (i, j) for any pair $(h, k) \in D$ can be different from 0 only if edge $e = [i, j]$ belongs to the circuit X . The objective function to be minimized is the sum of the design costs of the edges in X and of the routing costs of the arcs traversed by positive flow.

2.1 Multicommodity flow based formulation

The most natural way to formulate the problem is to use an integer multicommodity flow based model. For each edge $e = [i, j] \in E$ we introduce a binary variable x_e with value 1 if the edge $[i, j]$ is in the cycle X and 0 otherwise. We have then multicommodity flow variables for each commodity $(h, k) \in D$ and for each arc $(i, j) \in A$: f_{ij}^{hk} assumes value 1 if demand from h to k travels on arc (i, j) and 0 otherwise.

In order to simplify the formulation, we will use the following notation: given a subset of nodes $S \subset N$, $\delta(S) \subset E$ denotes the set of edges having exactly one endpoint in S and $\delta(i)$ is used as a shorthand of $\delta(\{i\})$. Similarly, $\delta^+(S) \subset A$ ($\delta^-(S) \subset A$) denotes the set of arcs having tail in S and head in $N \setminus S$ (head in S and tail in $N \setminus S$, respectively). We use $\delta^+(i)$ ($\delta^-(i)$) as a shorthand of $\delta^+(\{i\})$ ($\delta^-(\{i\})$). Given a subset $E^* \subseteq E$ and a subset $A^* \subseteq A$ we write $x(E^*)$ for $\sum_{e \in E^*} x_e$ and $f^{hk}(A^*)$ for $\sum_{(i,j) \in A^*} f_{ij}^{hk}$, respectively.

The TSP-GL is formulated as follows:

$$\min (1 - \alpha) \sum_{e \in E} \bar{c}_e x_e + \alpha \sum_{(i,j) \in A} \sum_{(h,k) \in D} q_{ij}^{hk} f_{ij}^{hk} \quad (1)$$

$$x(\delta(i)) = 2 \quad \forall i \in N \quad (2)$$

$$f^{hk}(\delta^+(i)) - f^{hk}(\delta^-(i)) = \begin{cases} 1 & \text{if } i = h, \\ -1 & \text{if } i = k, \\ 0 & \text{if } i \neq h, k \end{cases} \quad \forall (h, k) \in D, \forall i \in N \quad (3)$$

$$f_{ij}^{hk} \geq 0 \quad \forall (i, j) \in A, \forall (h, k) \in D \quad (4)$$

$$x_e - f_{ij}^{hk} \geq 0 \quad \forall e = [i, j] \in E, \forall (h, k) \in D \quad (5)$$

$$x_e - f_{ji}^{hk} \geq 0 \quad \forall e = [i, j] \in E, \forall (h, k) \in D \quad (6)$$

$$x_e \in \{0, 1\} \quad \forall e \in E \quad (7)$$

The objective function (1) accounts for the two components: the sum of design costs and the latency, where α is a suitable tradeoff parameter. Constraints (2) impose that the selected edges incident to each node must be exactly two, thus they must form a cycle cover of G . Constraint (3) are the multicommodity flow balance constraints stating that one unit of flow must be sent from h to k for every $(h, k) \in D$. Constraints (5) and (6) enforce the coherence between x and f variables, that is, any commodity (h, k) cannot travel on an arcs (i, j) or (j, i) if the corresponding edge $e = [i, j]$ is not in the cycle. Observe that constraints (5) and

(6) could be alternatively expressed in the apparently stronger form:

$$x_e - f_{ij}^{hk} - f_{ji}^{hk} \geq 0 \quad \forall e = [i, j] \in E, \forall (h, k) \in D. \quad (8)$$

There is no practical advantage in using this alternative expression however, because no optimal solution of the linear relaxation of (1)-(7) would have f_{ij}^{hk} and f_{ji}^{hk} strictly positive at the same time. We keep consequently the original formulation.

Note that, in general, our formulation does not avoid the presence of subtours. Subtours are made infeasible by the combination of coherence constraints and flow balance constraints if D is dense enough. The only case when the presence of subtours is possible is when D defines in G more than one connected component. In this case, in order to generate Hamiltonian circuits we must introduce explicit cutset inequalities. Usually, the demand matrix D in transit line design is dense enough to avoid this occurrence. In our investigation, we focus on the case of matrices D inducing a single connected component.

2.2 Related works

The TSP-GL can be seen as a specialized version of the FNDP consisting in finding a set of links (and capacities) on which flow demands can be routed, minimizing design and/or routing cost. In the case of the TSP-GL, links have infinite capacities and the network topology is constrained to be a Hamiltonian circuit. The FNDP is a very general framework able to model a wide variety of problems and for which several exact methods have been proposed. To cite a few: Benders decomposition (see Costa (2005) for a complete review on the topic), Lagrangean based techniques (see for example Crainic et al. 2001; Holmberg and Hellstrand 1998; Holmberg and Yuan 2000; Kliwer and Timajev 2005; Sellmann et al. 2002) and polyhedral approaches (see for example Atamtürk 2002; Barahona 1996; Bienstock and Günlük 1994; Günlük 1999; Leung and Magnanti 1989; Ortega and A. 2003; Raack et al. 2011).

More specifically, the Uncapacitated FNDP (UFNDP) has been addressed in Magnanti et al. (1986) by Benders decomposition. The authors applied a preprocessing phase based on a dual-ascent procedure that allowed to eliminate design variables and showed that the selection of Benders cuts is a key element for the success of the approach. Holmberg and Hellstrand (1998) proposed a Lagrangean relaxation of the flow conservation constraints and embedded such a relaxation scheme into a Branch and Bound algorithm. Ortega and A. (2003) considered a particular case of the UFNDP where one origin only is allowed and proposed a Branch and Cut algorithm. The so-called *dicut* inequalities showed to be very effective in their approach.

The TSP-GL also recalls a similar problem called Generalized Minimum Latency Problem introduced in Errico (2008). In the Generalized Minimum Latency Problem, we are given a complete directed graph $G = (N, A)$ where directed arcs are associated with a design cost and a routing cost. As in TSP-GL, origin/destination demand coefficients d_{hk} are given. The problem consists in determining a Hamiltonian directed cycle and in routing the demand on the cycle so that the sum of the design costs of the selected arcs and the routing costs multiplied by the flow is minimized. The main difference with respect to TSP-GL is that the cycle is directed, design costs are not necessarily symmetric, and demand routing must follow the direction of the arcs of the cycle.

Two important problems studied in the literature can be seen as special cases of the Generalized Minimum Latency Problem depending on how the demand is specified. One is the Minimum Latency Problem (MLP), also known in literature as the Traveling Deliveryman Problem or Traveling Repairman problem. A repairman has to visit a given set of customers, who issued a request before the repairman departure. Starting from a depot, the repairman has to visit all customers minimizing the total customer waiting time that depends on the distance traveled by the repairman before he/she serves them. The MLP can be modeled as a particular

case of Generalized Minimum Latency Problem where, supposing that the depot is in node 1, $d_{1k} = 1$ for $k = 2, \dots, n$ and 0 for all the other origin/destination pairs. Moreover the design costs are equal to zero for all arcs. The other problem is the Multicommodity TSP which is very similar to the MLP problem with the exception that the demand coefficients between node 1 and the others can assume any value greater than 0 rather than being unitary.

Lucena (1990) formulates the MLP using a time dependent TSP and obtains lower and upper bounds via dynamic programming and Lagrangean relaxation for instances with up to 30 nodes. More recently, Méndez-Díaz et al. (2008) proposed an evolution of the previous approach exploiting the polyhedral structure of the problem and increasing the size of the solved instances to 40 nodes. Different lower and upper bounds obtained using cumulative matroids and Lagrangean decomposition was proposed by Fischetti et al. (1993). The Multicommodity TSP is studied in Sarubbi and Luna (2007) and Sarubbi et al. (2008) where a classical Benders decomposition and a Quadratic Assignment approach are proposed.

We observe that Benders decomposition has been adopted for several problems related to the TSP-GL and this motivated our study. However, as detailed in Section 4, the particular structure of the TSP-GL does not allow for a straightforward extension of the known methods. We will exploit such a structure to develop an efficient cutting plane algorithm based on Benders decomposition and polyhedral properties of the TSP-GL, and embed it in an implicit enumeration scheme.

3 Benders decomposition and TSP-GL

The Benders decomposition was introduced in Benders (1962) as a method to solve mixed integer programs. It applies in particular to structured problems where we can identify a subset of “complicating” variables for which, when fixed, the problem reduces to an easy one. In our case it is easy to notice that if x variables defining the Hamiltonian circuit are fixed, the resulting restricted problem reduces to a set of independent minimum cost flow problems without capacity constraints. Let us analyze the Benders reformulation for the TSP-GL and illustrate the classical Benders decomposition.

Let us consider the following problem depending on x variables only:

$$\min (1 - \alpha) \sum_{e \in E} \bar{c}_e x_e + \alpha z(x) \quad (9)$$

$$x(\delta(i)) = 2 \quad \forall i \in N \quad (10)$$

$$x_e \in \{0, 1\} \quad \forall e \in E \quad (11)$$

Function $z(x)$ gives the contribution of the flow variables (i.e., the routing cost component) to the overall objective function and depends on the circuit definition. Given a possible solution \bar{x} , the value of $z(\bar{x})$ is obtained by solving the so called *PrimalSubproblem* (PS). For the TSP-GL, since there are no capacities on the arcs, the PS may be decomposed into independent one-commodity minimum cost flows, one for each origin-destination pair: $z(\bar{x}) = \sum_{(h,k) \in D} z^{hk}(\bar{x})$, where

$$z^{hk}(\bar{x}) = \min \sum_{(i,j) \in A} q_{ij}^{hk} f_{ij}^{hk} \quad (12)$$

$$f^{hk}(\delta^+(i)) - f^{hk}(\delta^-(i)) = \begin{cases} 1 & i = h, \\ -1 & i = k, \\ 0 & \forall i \in N \setminus \{h, k\} \end{cases} \quad \forall i \in N \quad (13)$$

$$-f_{ij}^{hk} \geq -\bar{x}_e \quad \forall i, j \in N, i < j, e = [i, j] \quad (14)$$

$$-f_{ji}^{hk} \geq -\bar{x}_e \quad \forall i, j \in N, i < j, e = [i, j] \quad (15)$$

$$f_{ij}^{hk} \geq 0 \quad \forall (i, j) \in A. \quad (16)$$

Observe that the strong duality holds for every flow subproblem, hence, we can equivalently consider the *Dual Subproblems* :

$$z^{hk}(\bar{x}) = \max \quad \rho_h^{hk} - \rho_k^{hk} - \sum_{e \in E} \bar{\lambda}_e^{hk} \bar{x}_e \quad (17)$$

$$\rho_i^{hk} - \rho_j^{hk} - \lambda_{ij}^{hk} \leq q_{ij}^{hk}, \quad \forall (i, j) \in A \quad (18)$$

$$\lambda_{ij}^{hk} \geq 0 \quad \forall (i, j) \in A \quad (19)$$

$$\rho_i^{hk} \geq 0 \quad \forall (i, j) \in A \quad (20)$$

where $\bar{\lambda}_e := \lambda_{ij} + \lambda_{ji}$, $\forall e = [i, j] \in E$.

By denoting $Q = \{(\lambda, \rho) \in \mathcal{R}^{(|A|+|N|)|D|} \mid (\lambda, \rho) \text{ satisfy (17), \dots, (20)} \forall (h, k) \in D\}$, problem (9) - (11) can be rewritten as:

$$\min_x (1 - \alpha) \sum_{e \in E} \bar{c}_e x_e + \alpha \max_{(\lambda, \rho) \in Q} \left\{ \sum_{(h, k) \in D} (\rho_h^{hk} - \rho_k^{hk} - \sum_{e \in E} \bar{\lambda}_e^{hk} x_e) \right\} \quad (21)$$

$$x(\delta(i)) = 2 \quad \forall i \in N \quad (22)$$

$$x_e \in \{0, 1\} \quad \forall e \in E \quad (23)$$

The fact that flow variables f do not explicitly appear in problem in (9) - (11), and as a consequence also in problem (21) - (23), implies that coherence constraints linking x and f can be violated. Solutions x that are feasible for (9) - (11) might therefore not guarantee the feasibility of the PS, and DS would result to be unbounded. To avoid this circumstance, we impose in problem (9) - (11) (and equivalently in problem (21) - (23)) that $\sum_{(h, k) \in D} (\rho_h^{hk} - \rho_k^{hk} - \sum_{e \in E} \bar{\lambda}_e^{hk} x_e) \leq 0$ for every (λ, ρ) extreme ray of Q .

Since unbounded solutions are now avoided, the new DS has optimum solution in one of the extreme points of Q . In order to obtain the Benders reformulation, we linearize the objective function (21). Because the number of extreme points of a polyhedron is finite, the expression (21) can be seen as the minimum of a set of affine functions which could be linearized with the addition of a continuous variable η .

Now we can write the Benders reformulation of the TSP-GL:

$$\min \quad \alpha \eta + (1 - \alpha) \sum_{e \in E} \bar{c}_e x_e \quad (24)$$

$$\eta + \sum_{e \in E} \sum_{(h, k) \in D} \bar{\lambda}_e^{hk} x_e \geq \sum_{(h, k) \in D} (\rho_h^{hk} - \rho_k^{hk}) \quad \forall (\lambda, \rho) \in \text{extr}(Q) \quad (25)$$

$$\sum_{e \in E} \sum_{(h, k) \in D} \bar{\lambda}_e^{hk} x_e \geq \sum_{(h, k) \in D} (\rho_h^{hk} - \rho_k^{hk}) \quad \forall (\lambda, \rho) \in \text{rays}(Q) \quad (26)$$

$$x(\delta(i)) = 2 \quad \forall i \in N \quad (27)$$

$$x_e \in \{0, 1\} \quad \forall e \in E \quad (28)$$

where $\text{extr}(Q)$ and $\text{rays}(Q)$ are the set of extreme points and rays of Q respectively.

Notice that the problem (24)-(28) does not contain f variables, but the number of constraints (25) and (26) (called *optimality* and *feasibility* cuts, respectively) is very high and this suggests the use of a cutting plane

approach. In fact Benders decomposition is nothing but a cutting plane algorithm applied to the Benders reformulation. The classic version of the algorithm considers a relaxation of the Benders reformulation, called *Master Problem* (MP), with only a small subset of optimality and feasibility cuts. The master problem, which is a MIP itself, is solved. If the resulting DS is unbounded, then a feasibility cut is added. If DS is bounded, either the current solution is optimal or it is possible to identify a violated optimality cut. The process is then iterated.

4 Three Benders-based algorithms

We now describe three algorithms addressing the TSP-GL based on Benders decomposition. In Section 4.1 we present the main motivations and formalize the common algorithmic scheme, called Benders-Branch-and-Cut (BBC). The three versions of the BBC differ on the type of generated feasibility cuts. In Section 4.2, we describe a basic version of the algorithm (BBC1) based on the classical Benders reformulation (24)-(28) and prove some mathematical properties. In Section 4.3, we underline the common points between TSP-GL and S-TSP and describe two algorithms that take advantage of such similarities (BBC2 and BBC3).

4.1 General algorithmic approach and motivations

The classic version of the Benders decomposition (CBD) requires the solution of a MIP at *every* master iteration. The drawback is that the complexity of such MIPs rapidly increases with the number of added constraints, which usually spoil the structure of the original problem, and convergence to optimality soon becomes impractical, except for small instances. For this reason, research focused on obtaining Benders cuts in alternative ways so as to reach a better convergence. One of the most popular techniques was proposed by McDaniel and Devine (1977). Starting from the observation that the feasible region of a dual Benders subproblem does not depend on the particular master problem solution, the authors deduce that any extreme point or ray of a dual subproblem, independently on how it was obtained, can be used to generate valid Benders cuts. Such a property implies, for example, that the solution obtained by solving the *linear relaxation* of the MP can be used to obtain valid Benders cuts, with the advantage of being much easier to solve than the corresponding integer version. McDaniel and Devine (1977) proposed a modification of the Benders decomposition where the initial cuts were obtained by applying the Benders scheme to the linear relaxation of the MP. When certain conditions apply, the authors propose to switch to CBD. Let us call MBD this type of solution approach. MBD has been successfully applied in several works, as for example in Cordeau et al. (2000, 2001); Rei et al. (2009).

Observe that stopping the continuous phase of the MBD when no Benders cuts can be found, corresponds to solving the linear relaxation of the initial formulation via Benders decomposition. It is actually possible to further modify MBD to progressively recover integrality by embedding the linear relaxation of the MP in a branch-and-bound scheme, instead of switching to the CBD. We name such a solution scheme Benders-Branch-and-Cut algorithm (BBC). BBC has been adopted in several works, as for example in Ljubić et al. (2012); Rodríguez-Martín and Salazar-González (2008); Sridhar and Park (2000). In Naoum-Sawaya and Elhedhli (2012), BBC showed to be up to seven time faster than CBD with Pareto-optimal cuts (for details on Pareto-optimal cuts, see Magnanti and Wong 1981) when applied to capacitated facility location problem. Our preliminary attempts to solve TSP-GL by CBD and MBD could only solve very small instances of TSP-GL, even if a certain number of additional algorithmic improvements were implemented (Pareto-optimal cuts, heuristic generation of Benders cuts, warm-restart, etc).

For the above reasons, the algorithms that we present here follow the BBC scheme. In order to handle

the enumeration, the BBC algorithm makes use of a list of candidate problems to be solved and can be summarized as follows:

- STEP 0 Set iteration counter $\mu = 0$, global upper bound $UB = +\infty$. In the initial MP integrality constraints are relaxed and only degree constraints (2) are considered. In addition, a set of optimality cuts obtained heuristically is added to the formulation. The resulting starting MP is inserted into the list of problems.
- STEP 1 If problem list is empty, exit. Else, select an MP from the list. Set the current problem upper bound $problemUB = UB$.
- STEP 2 Set $\mu = \mu + 1$. Solve MP and let (η^μ, x^μ) be its solution.
- STEP 3 Execute the feasibility cut separation procedure. If x^μ violates some feasibility cut add it to MP and go back to STEP 2. Else go to STEP 4;
- STEP 4 Execute the optimality separation procedure solving the Benders subproblem. Let f^μ be the corresponding flow. If $(qf^\mu + cx^\mu < problemUB)$, update $problemUB$. If $(\eta^\mu + cx^\mu < nodeUB - \epsilon)$, a violated optimality cut has been found, add it to MP and go back to STEP 2. Else, current MP optimality is reached, go to STEP 5.
- STEP 5 If $(\eta^\mu + cx^\mu > UB - \epsilon)$ fathom the current node. Else, If (f^μ, x^μ) is integer, update UB . If it is not integer, a branching procedure is performed and the new problems are added to the list. Go to STEP 1.

4.2 BBC1: The basic version

The BBC1 algorithm implements the BBC scheme and adopts feasibility cuts defined by (26). Let us denote z_1 the value of the optimal solution of the linear relaxation of the MP in BBC1, and z_{MF} the value of the linear relaxation of the multicommodity flow formulation (1)-(7). It is known that the classical Benders reformulation (24)-(28) is equivalent to the original problem in terms of lower bounds provided by their linear relaxation, hence $z_1 = z_{MF}$.

The natural extension to the full space of any extreme ray of $Q^{hk} = \{(\lambda^{hk}, \rho^{hk}) \in \mathcal{R}^{|A|+|N|} \mid (\lambda^{hk}, \rho^{hk}) \text{ satisfy (18), \dots, (20)}\}$ for any given $(h, k) \in D$ is a ray of Q itself (for a formal proof about sets of ray generators of cones with diagonal block structure see Padberg and Sung 1991). Moreover, recall that the Benders subproblem decomposes into smaller independent subproblems, one per commodity. Constraints (26) can be consequently substituted by the following disaggregated version:

$$\rho_h^{hk} - \rho_k^{hk} - \sum_{e \in E} \bar{\lambda}_e^{hk} x_e \leq 0 \quad \forall (\lambda^{hk}, \rho^{hk}) \in \text{extr}(Q^{hk}), \forall (h, k) \in D. \quad (29)$$

The following proposition provides a result regarding the strength of these constraints.

Proposition 4.1. *Let us denote z_{CI} the value of the optimal linear relaxation of (24), (25), (27), (28) and the following Cutset Inequalities:*

$$x(\delta(S)) \geq 1 \quad \forall S \subset V. \quad (30)$$

Then, $z_1 = z_{CI}$

Proof. Proof. We prove the proposition by showing that inequalities (29) are equivalent to (30). To this purpose, we recall that Padberg and Sung (1991) provide a characterization of the extreme rays of the node-arc cone of a directed graph, i.e., the extreme rays of Q^{hk} . In particular, they prove that the extreme rays of Q^{hk} are positive multiple of vectors $(\lambda^{hk}, \rho^{hk})$ given by (i), (ii), (iii) where

$$\begin{aligned}
(i) \quad \rho_i^{hk} &= 0 \quad \forall i \in N, & \lambda_{ij}^{hk} &= \begin{cases} 1 & \text{for exactly one } i \text{ and one } j \in N \\ 0 & \text{otherwise} \end{cases} \\
(ii) \quad \rho_i^{hk} &= \begin{cases} 1 & \text{for all } i \in S \\ 0 & \text{otherwise} \end{cases}, & \lambda_{ij}^{hk} &= \begin{cases} 1 & \text{for all } i \in N \setminus S, j \in S \\ 0 & \text{otherwise} \end{cases} \\
(iii) \quad \rho_i^{hk} &= \begin{cases} -1 & \text{for all } i \in S \\ 0 & \text{otherwise} \end{cases}, & \lambda_{ij}^{hk} &= \begin{cases} 1 & \text{for all } i \in S, j \in N \setminus S \\ 0 & \text{otherwise} \end{cases}
\end{aligned}$$

and $S \subseteq N, 1 \leq |S| \leq n - 1$. The substitution of (i) in (29) gives the redundant inequality $x_e \geq 0, \forall e \in E$. Let us consider now extreme rays of the form (ii) and the case $h \in S, k \in N \setminus S$. Substituting in (29) one obtains the desired inequality. If both h and k belong to S or to $N \setminus S$ we obtain redundant equations. The case (iii) is similar to (ii). Because we assumed that the demand matrix D induces a single connected component, the proposition is proved. \square

From the algorithmic point of view, BBC1 separates inequalities (29) by solving a minimum cost flow algorithm for each commodity, using a network simplex algorithm, for example. As soon as a violated inequality is found, the separation algorithm is stopped and the corresponding cut is added.

4.3 BBC2 and BBC3: Exploiting the Hamiltonian circuit structure

The formulation (1)-(7) does not explicitly exploit the particular Hamiltonian circuit structure of the design part of TSP-GL. An important property that plays a central role for BBC2 and BBC3 is pointed out by the following proposition.

Proposition 4.2. *Every valid inequality for the S-TSP is also valid for the TSP-GL.*

Proof. Proof. Let P denote the polyhedron defined by constraints (2)–(7), and let P_{S-TSP} be the polyhedron of the Hamiltonian circuits defined on the same graph G . Assuming that the demand matrix D is dense enough to make subtours infeasible, calling $Proj_x(P) = \{x \in \{0, 1\}^{|E|} \mid \exists f \in \mathcal{R}^{|A||D|}, (x, f) \in P\}$ we have that $Proj_x(P) \subseteq P_{S-TSP}$. Moreover, given any $x \in P_{S-TSP}$ it is always possible to find a flow $f \in \mathcal{R}^{|A|^2}$ such that the corresponding (x, f) belongs to P . As a consequence, $P_{S-TSP} \equiv Proj_x(P)$ and the fact that every valid inequality for the S-TSP is also valid for the TSP-GL follows immediately. \square

There has been an extensive research effort to study the S-TSP polyhedron and a number of families of facet defining inequalities have been found. The above result suggests to exploit such knowledge about the S-TSP polyhedron to possibly strengthen the lower bound given by the linear relaxation of the MP. To this purpose, consider the well-known family of S-TSP facet defining inequalities called Subtour Elimination Constraints (SEC):

$$x(\delta(S)) \geq 2 \quad \forall S \subset V, \tag{31}$$

and let us call P_{MF} the polyhedron of the linear relaxation (1)-(7), and P_{SEC} the polyhedron defined by P_{MF} with the addition of inequalities (31). Obviously, $P_{SEC} \subseteq P_{MF}$, however the inclusion is tight, as shown by the following proposition.

Proposition 4.3. *P_{SEC} is strictly included in P_{MF}*

Proof. Proof. Let $Proj_x(P_{MF}) = \{x \in \{0,1\}^{|E|} \mid \exists f \in \mathcal{R}^{|A||D|}, (x, f) \in P_{MF}\}$ and similarly $Proj_x(P_{SEC}) = \{x \in \{0,1\}^{|E|} \mid \exists f \in \mathcal{R}^{|A||D|}, (x, f) \in P_{SEC}\}$. We prove the proposition by showing that $Proj_x(P_{SEC})$ is strictly included in $Proj_x(P_{MF})$. Recall that Benders reformulation allows us to write $Proj_x(P_{MF}) = \{x \in \{0,1\}^{|E|} \mid (26) \text{ and } (27) \text{ are satisfied}\}$ and in Proposition 4.1 we saw that this is equivalent to $Proj_x(P_{MF}) = \{x \in \{0,1\}^{|E|} \mid (30) \text{ and } (27) \text{ are satisfied}\}$. By definition, $Proj_x(P_{SEC}) = \{x \in \{0,1\}^{|E|} \mid (26), (27) \text{ and } (31) \text{ are satisfied}\}$. The fact that inequalities (30) are clearly implied by (31) together with the observation that it is trivial to find a solution belonging to $Proj_x(P_{MF})$ and not to $Proj_x(P_{SEC})$, proves that $Proj_x(P_{SEC})$ is strictly included in $Proj_x(P_{LP})$. \square

Proposition 4.3 implies that inequalities (31) can be exploited to improve the lower bound provided by BBC1.

Let us denote by P_{FAC} the polyhedron defined by including an additional set of families of S-TSP facet defining inequalities into P_{SEC} . Because the added inequalities are facet defining for the S-TSP and because in Proposition 4.2 we saw that the projection on the x -space of P_{MF} is the S-TSP polyhedron, the following proposition easily follows:

Proposition 4.4. *P_{FAC} is strictly included in P_{SEC} .*

Therefore we may conclude that the following hierarchy of bounds holds true:

Corollary 4.1. *Denote z_{SEC} and z_{FAC} the optimal value of minimizing the objective function (1) over polyhedra P_{SEC} and P_{FAC} , respectively. Then, $z_1 = z_{MF} \leq z_{SEC} \leq z_{FAC}$.*

However, note that the improvement of z_{MF} by adding S-TSP cuts to the original formulation (1)-(7) is far from being practical. In fact, as we will point out in Section 5.2, solving the linear relaxation of the original formulation is extremely time consuming and solution techniques requiring the iterated solution of such a problem cannot be efficient. In this sense, Benders reformulation is a key element of the success of our approach because it projects out the flow variables and makes it possible, in a BBC framework, to exploit cutting plane generation and in particular the Hamiltonian structure of the network design in the TSP-GL.

BBC2 and BBC3 implement the general BBC algorithm and differ in the families of S-TSP cuts generated. In the following, we briefly list inequalities and the adopted separation algorithms.

4.3.1 BBC2.

BBC2 is similar to BBC1, but adopts inequalities (31) as feasibility cuts instead of (26). Denoting z_2 the optimal value of the linear relaxation provided by algorithm BBC2, it is clear that $z_2 = z_{SEC}$.

For the separation of inequalities (31), BBC2 adopts the exact method described in Padberg and Rinaldi (1990).

4.3.2 BBC3.

In addition to constraints (31) introduced by BBC2, BBC3 considers a variety of additional families of S-TSP facet defining inequalities as a subset of comb inequalities and local cuts. Let us denote z_3 the optimal value of the linear relaxation provided by algorithm BBC3. From Corollary 4.1, we can clearly deduce that $z_3 \geq z_2$.

Recall that a *comb* C is a subset of N and is specified by a *handle* $H \subseteq N$ and an odd number of disjoint *teeth* $S_1, \dots, S_k \subseteq N$ such that each tooth has at least one node in H and at least one node not in H . For every comb C , the corresponding inequality is:

$$x(\delta(H)) + \sum_{j=1}^k x(\delta(T_j)) \geq 2k + 1. \quad (32)$$

It is not known whether the general form of comb inequalities admits a polynomial-time algorithm. However, polynomial algorithms or fast heuristics for specific comb inequalities have been proposed. BBC3 implements both a heuristic and an exact algorithm described in Padberg and Rao (1982) for *blossoms*, i.e., combs whose teeth have only two nodes. BBC3 also implements a heuristic separation of comb inequalities starting from *blocks*, i.e., maximal two-connected components of the residual graph induced by the relaxed solution. Comb inequalities from blocks are separated following the algorithm in Naddef and Thienel (2002). Finally, BBC3 implements separation procedures for *local* cuts. Such cuts do not follow a predefined template (as subtour elimination or comb inequalities), but are obtained by suitably shrinking nodes of the residual graph induced by a fractional solution and mapping violated inequalities back in the original graph representation. Details on local cuts can be found in Applegate et al. (2007), Chapter 11.

5 Computational experience

We perform three experimental campaigns with the scope of 1) comparing lower bounds and their computational efficiency, 2) selecting the overall best algorithm, 3) testing the selected algorithm on suitable benchmark instances. More specifically, the first campaign aims at choosing which algorithm among BBC1, BBC2, and BBC3 gives the best performance in terms of root bounds and computing times. For completeness, we also compare BBC algorithms with a state-of-the-art LP solver applied to (1)-(7). Since there is computational evidence of the superiority of BBC3, in the second campaign we test several possible implementation options applied to BBC3 with the scope of identifying a unique algorithmic setting that averagely outperforms the others. We call the resulting algorithm BBC3+. In the third experimental campaign, we test BBC3+ on a new set of benchmark instances suitably derived from the well-known TSPLIB library. Results show that benchmark instances with up to 42 nodes are solved to optimality, and, in the worst case, a feasible solution within 8.2% of the optimal value is obtained on a 101-nodes instance.

Given the different scopes and natures of the three experimental campaign, a suitable set of instances has been chosen for each campaign. However, construction criteria are mostly the same for all sets, and inspired by two main principles: 1) to be as general as possible 2) be representative of our main application context, i.e., the design of semi-flexible transit lines.

5.1 Instance sets

We built two instance sets each with a specific purpose: TSPGL1 to compare algorithms and algorithmic options. To guarantee statistical robustness, TSPGL1 includes a wide variety of problem types, and for each type, several randomly generated instances are considered. TSPGL2 provides the first benchmark instance

set for the TSP-GL. For this, TSPGL2 builds upon the widely used benchmark TSPLIB and includes one instance only for each problem type. For both instance sets, the problem dimension range reflects typical situations found in our reference application.

Each instance in TSPGL1 is characterized by a set of nodes and a demand matrix D . The nodes are uniformly generated in a square of side length 100 and the graph is complete. c_{ij} is given by the Euclidean distance between nodes i and j , while c_e is set equal to twice the distance between the two endpoints. This choice comes from the reference application of the TSP-GL, i.e., the design of semi-flexible transit lines, (see Errico et al. 2011b, for details), where two vehicles are considered to travel the transit line in opposite directions. The routing costs q_{ij}^{hk} on the arcs, that model passenger travel times, are set to $c_{ij} \frac{d_{hk}}{\sum_{(s,t) \in D} d_{st}}$, where $\frac{d_{hk}}{\sum_{(s,t) \in D} d_{st}}$ is the contribution of demand (h,k) . Therefore $\sum_{(i,j) \in A} \sum_{(h,k) \in D} q_{ij}^{hk} f_{ij}^{hk}$ measures the average passenger travel time.

As for demand matrices, we considered two different classes, reflecting the typical application context in which a few nodes, called *poles*, concentrate most of the transit demand. In particular, we considered the following scenarios:

- **C**: Complete and Polarized: The demand matrix is complete and each d_{hk} follows a uniform distribution on a given interval. However, a few poles are randomly picked among the nodes and approximately half of the total demand volume has origin *and* destination among poles.
- **S**: Sparse and Polarized: As in C, but the demand matrix is sparse and the number of destinations associated with a given node is limited to 4 and chosen randomly.

As specified in Errico et al. (2011b) each of the above demand matrix classes corresponds to a different approach in the design of a semi-flexible transit line, where C and S account for more or less aggregated data, respectively.

Observe that, for completeness, we initially considered also the following two set of demand matrices:

- Complete and Uniform: Demand matrix is complete and d_{hk} are generated uniformly on a given interval.
- Sparse and Uniform: As above, but the demand matrix is sparse as the number of destinations associated with a given node is limited to 4 and chosen randomly.

The results obtained on these instances were not qualitatively different from the corresponding polarized ones.

TSPGL1 considers networks with the following number of nodes: from 15 to 40 included, by steps of 5, and 50, 60, 80, and 100. We also considered three possible values of α in the objective function: 0.25 (emphasis on design), 0.75 (emphasis on latency), 0.5 (fairness). Finally, for each network dimension, demand setting, and objective emphasis, we generated 10 instances considering different random seeds for the generation of random data. As a result, TSPGL1 is composed of 720 different instances.

The building process for TSPGL2 is quite similar to TSPGL1, with a few important differences. First of all, the networks are deduced from the TSPLIB library instead of being generated randomly. We selected all the TSPLIB instances of the S-TSP with at most 101 cities. As in the case of TSPGL1, we considered two demand scenarios for each network, C and S. Differently from TSPGL1, TSPGL2 only considers $\alpha = 0.5$ in the objective function and one only representative instance is considered for each network/demand/emphasis combination. The resulting 48 instances are available online (TSPGL2 2012).

5.2 Comparing lower bounds and efficiency

In this section, we compare root lower bounds and efficiency provided by algorithms BBC1, BBC2 and BBC3. For completeness, we considered also the solution of the linear relaxation of the original formulation (1)-(7) by a commercial solver, more specifically, we used the dual simplex solver in CPLEX 12.3.0.0, which gave better performances with respect to other algorithms. In the following, MF will denote these results.

We implemented the BBC algorithm in the C++ programming language. The master problem was solved by the MIP solver available in CPLEX 12.3.0.0 with Concert Technology. The dual simplex was used as linear relaxation solver. Separation procedures have been implemented as instantiations of several user-callback available in Concert Technology. In particular, for BBC1 the Benders subproblem has been solved via Network Simplex Optimizer implemented in CPLEX 12.0.0.3. S-TSP separation procedures in BBC2 and BBC3, i.e., subtour elimination, heuristic and exact blossom, combs from blocks, and local cuts have all been implemented by suitably importing in our codes separation procedures of the well-known S-TSP solver Concorde, version 03-12-19, coded in C programming language (Applegate et al. 2003). Details on the latter set of separation procedures can be found in Applegate et al. (2007). Experiments were executed on a Dual-Core AMD Opteron(tm) Processor 2218, with 4G DDR2 RAM memory on SunOS 5.10.

Instance			GAP			Times (s)			
dimension	D. type	α	BBC1	BBC2	BBC3	CPLEX	BBC1	BBC2	BBC3
15	C	0.25	3.4	73.6	78.3	5.4	0.9	0.9	0.9
20	C	0.25	3.2	82.3	90.6	40.1	5.0	3.3	4.5
25	C	0.25	3.3	69.1	85.3	265.9	17.1	9.6	12.9
30	C	0.25	3.5	62.8	80.3	1309.0	46.2	28.4	44.1
35	C	0.25	3.5	68.2	89.0	2872.3	75.8	50.5	57.4
15	C	0.5	2.8	66.8	69.1	8.4	2.4	1.8	2.4
20	C	0.5	3.5	52.5	61.1	69.4	13.4	7.9	11.7
25	C	0.5	4.3	45.3	57.8	486.9	39.7	20.2	36.3
30	C	0.5	4.9	33.7	50.4	1827.6	115.1	69.4	101.5
35	C	0.5	4.4	41.7	62.7	8779.2	277.8	157.8	178.7
15	C	0.75	3.4	25.9	29.6	16.2	8.5	5.4	11.4
20	C	0.75	6.6	20.9	26.0	98.7	32.3	26.5	49.4
25	C	0.75	7.8	15.4	23.3	862.6	109.4	77.7	109.1
30	C	0.75	9.2	13.2	19.7	4749.9	340.2	245.6	301.8
35	C	0.75	8.0	17.2	23.8	17495.8	742.8	494.9	470.2
15	S	0.25	3.4	74.4	77.2	5.3	0.6	0.5	0.5
20	S	0.25	3.2	80.0	89.9	22.1	2.7	1.8	2.5
25	S	0.25	3.4	67.9	84.0	149.2	7.7	4.3	6.6
30	S	0.25	3.5	62.7	79.8	552.4	22.2	12.7	18.9
35	S	0.25	3.5	67.9	88.0	1318.1	30.9	23.8	24.3
15	S	0.5	2.9	67.4	65.9	6.8	1.6	1.1	1.6
20	S	0.5	3.7	53.5	61.7	38.7	6.5	4.3	6.6
25	S	0.5	4.4	44.2	56.2	203.3	19.5	10.2	18.4
30	S	0.5	5.0	34.6	50.6	738.4	48.4	31.4	48.9
35	S	0.5	4.6	41.6	60.9	2339.8	90.6	54.0	77.1
15	S	0.75	3.5	28.1	31.6	11.8	5.2	3.1	7.9
20	S	0.75	6.9	23.3	28.1	50.9	16.6	14.7	33.5
25	S	0.75	8.5	16.3	22.9	276.2	51.4	33.7	70.9
30	S	0.75	9.5	13.4	19.7	1306.0	145.0	100.0	160.9
35	S	0.75	8.4	17.4	23.9	3397.3	275.7	189.5	270.3

Table 1: Comparison among CPLEX and BBC1, BBC2, and BBC3 on small instances

We report in Table 1 the numerical results about the comparison of MF with BBC1, BBC2, and BBC3 on small instances of the TSPGL1 set. Recalling that TSPGL1 includes 10 random samples for each instance type, rows in Table 1 report average results for each instance type. In the first three columns the instance type is identified by indicating number of nodes, demand type, and value of α . The following column reports the relative percentage gap between the value of the optimal solution and the lower bound of the linear relaxation provided by z_{MF} . Given that, as observed in Section 4.2, $z_{MF} = z_1$, there is no need to report

the lower bound quality of BBC1. The next two columns display the average gap percentage closed by z_2 and z_3 with respect to z_{MF} . The last four columns report the average computing times in seconds.

We observe that the lower bound provided by BBC1 (equivalent to MF) is the worst and has a gap ranging from about 3% to no more than 10%. This gap increases with size of the instances and is also affected by the value of α . When α is higher, i.e., when the travel time component of the objective function is greater, the gap increases, while when the design component dominates, thus when the instances are closer to a TSP problem, the bound is tighter. Note also that the gap of BBC1 is not affected by the type (C/S) of the demand matrix. As for the computing times, there is a huge difference between the commercial software and BBC1. CPLEX times dramatically increase with the size and also testify that instances with high values of α are more difficult. Moreover, there is a great impact on the times due to the demand matrix type. Dense instances (C) require much higher times with respect to the sparse ones (S). This is mainly due to the number of variables of the model. Such a difference in computing times is mitigated using BBC1 instead of CPLEX. Besides requiring much smaller times, BBC1 reveals also to be more stable with respect all parameters.

Considering the other two bounds, we may observe a clear improvement in terms of gap and computing time. The gap closed by BBC2 with respect to BBC1 ranges in average from 80% for the easy instances to 13% for the more difficult ones. Notice also that the time of BBC2 is shorter than that of BBC1. BBC3 reaches better bounds closing a larger amount of gap and has computing times comparable with BBC1 and sometimes smaller for big instances, thus showing a better trend with respect to the instance size.

It is interesting to notice how the TSP-GL behaves differently than the general FNDP family. In the latter problems in fact, the more important the flow component with respect to the design component is, the easiest is the problem. This is normally explained noticing that the linear part of the problem (the flows) becomes predominant with respect to the non linear one (the design). For the TSP-GL the situation is different because the Hamiltonian structure of the network topology seems to emphasize the non-linearity of the problem when flow costs are higher. From the computational view point, we observed for BBC algorithms that the number of optimality cuts generated is much higher for high values of α . This suggests that the optimal solution of the master problem is located towards the interior of the feasible region of the projection in the x -subspace of the TSP-GL polyhedron. This would also explain why optimal solutions of the linear relaxation are more fractional, with consequently worse lower bounds, and why S-TSP facet defining inequalities are less helpful for high α .

We performed a second experimental test where we compared BBC1, BBC2 and BBC2 on instances of TSPGL1 with a higher number of nodes, and limiting the computing time to 5 hours. In this second test, we did not compare with MF because the commercial solver usually failed due to memory limits. Also, differently from the previous experimentation, we do not know the values of the optimal solutions of the integer problems. We consequently compared the relative improvements of lower bound values instead of comparing the relative gaps with the optimum. Results are reported in Table 2.

The meaning of the columns is as follows: the first three columns identify the instance type by specifying the number of nodes, the demand type and the value of α , respectively. The next column reports the average value of BBC1, on the 10 considered random seeds, while the following two columns report the average relative improvement of BBC2 and BBC3 over BBC1, respectively, computed as $100(z_i - z_1)/z_1$, $i = 2, 3$. The next three columns represent the average computing times, while the last three columns report the number instances for which the algorithm stops due to time limit reasons.

Table 2 confirms the trends observed in Table 1. Perhaps the most important fact underlined by Table 2 is that BBC3 consistently outperforms BBC1 and BBC2 in terms of lower bounds. Its computing times are comparable with BBC2, and are even better on harder instances, as it can be deduced by comparing how many times BBC3 hits the time limit with respect to BBC1 and BBC2.

Instance			GAP	GAP %Impr.			Times (sec)			#Time limit		
dimension	D. type	α	BBC1	BBC2	BBC3	BBC1	BBC2	BBC3	BBC1	BBC2	BBC3	
40	C	0.25	761.8	2.6	3.4	166	95	132	0	0	0	
50	C	0.25	848.1	3.3	4.1	599	351	469	0	0	0	
60	C	0.25	916.6	3.1	4.0	1354	1079	1177	0	0	0	
80	C	0.25	1040.8	3.1	4.0	7266	4951	6417	0	0	0	
100	C	0.25	1148.0	3.4	4.3	17937	16392	16219	9	7	6	
40	C	0.5	538.2	2.3	2.9	509	393	353	0	0	0	
50	C	0.5	597.0	2.9	3.6	2060	1505	1297	0	0	0	
60	C	0.5	640.7	2.6	3.4	6466	4105	3340	1	0	0	
80	C	0.5	726.7	2.5	3.3	17604	15836	15906	8	6	4	
100	C	0.5	793.7	3.3	4.1	18000	18000	18000	10	10	10	
40	C	0.75	309.3	1.7	2.1	1371	1088	1062	0	0	0	
50	C	0.75	339.9	2.2	2.6	5378	4805	4237	0	0	0	
60	C	0.75	358.9	1.9	2.4	15634	11760	10576	4	0	0	
80	C	0.75	401.8	2.5	3.2	18000	18000	18000	10	10	10	
100	C	0.75	414.0	5.8	7.0	18000	18000	18000	10	10	10	
40	S	0.25	761.3	2.7	3.5	1860	38	53	1	0	0	
50	S	0.25	847.6	3.3	4.1	203	116	164	0	0	0	
60	S	0.25	915.7	3.1	4.0	492	345	404	0	0	0	
80	S	0.25	1040.6	3.1	4.0	2477	1398	1632	0	0	0	
100	S	0.25	1147.4	3.4	4.3	8105	5896	6916	0	0	0	
40	S	0.5	537.3	2.4	3.0	1975	144	154	1	0	0	
50	S	0.5	596.0	3.0	3.6	637	484	484	0	0	0	
60	S	0.5	639.2	2.8	3.5	1525	1154	1269	0	0	0	
80	S	0.5	726.3	2.5	3.4	7561	4886	5026	0	0	0	
100	S	0.5	796.1	3.0	3.8	17932	17490	15872	9	9	5	
40	S	0.75	308.0	1.8	2.2	494	375	485	0	0	0	
50	S	0.75	338.3	2.3	2.7	1681	1448	1526	0	0	0	
60	S	0.75	356.9	2.1	2.6	6429	3905	3849	1	0	0	
80	S	0.75	405.1	1.8	2.5	17988	16162	16057	9	4	4	
100	S	0.75	433.3	3.2	3.9	18000	18000	18000	10	10	10	

Table 2: Comparison among BBC1, BBC2 and BBC3 on bigger INST1 instances. Time limit: 5h

Summarizing, we can conclude that BBC algorithms clearly outperform MF solved by a state-of-the-art linear programming solver. Overall, BBC3 favorably compares with BBC1 and BBC2 because it yields better lower bounds and competitive computing times. For this reason, the following experimental campaigns focus on BBC3 only.

5.3 Implementation refinements and details

The efficiency and effectiveness of Benders decomposition-based approaches strongly depend on the quality of the implementation. For this reason, we considered several options and compared their impact on the solution of our problem. This section aims to evaluate the impact of single algorithmic options or refinements on the performance of BBC3.

Given the high number of possible algorithmic combinations, we decided to perform the tests on a restricted but representative subset of TSPGL1. In particular, we only considered instances with 40, 50 and, 60 nodes. Instance types as well as the number of seeds used for each type are as in the previous Section.

The rest of the section is organized as follows: In Section 5.3.1, we discuss the generation of the initial sets of Benders cuts and heuristic strategies. In Section 5.3.2, we consider the generation of disaggregated optimality cuts (multi-cut). In Sections 5.3.3 and 5.3.4, we compare alternative separation and branching strategies.

5.3.1 Initial cuts and heuristics.

Heuristic methods might improve the convergence of the BBC algorithm for at least two reasons: 1) They can provide good solutions in early stages of the algorithm and this can reduce the size of the branch-and-bound tree. 2) Heuristic solutions can be used to generate optimality cuts in an easy way without having to solve STEP 2 of the BBC algorithm. This can be of some help both in the starting phase of the algorithm, to provide good initial cuts and starting solutions, and during the algorithm.

We implemented a simple local search based on the traditional two-opt exchange move for the S-TSP. Observe in fact that any Hamiltonian circuit can be easily extended to obtain a feasible TSP-GL solution. In particular, given a solution of the TSP-GL, we apply a two-opt procedure to the circuit defined by the x variables, update the optimal flows and evaluate the new costs. Such a procedure can be applied at several points of the algorithm. We implemented and compared the following variants of the BBC3 algorithm:

- **h1:** No initial solution or cut is provided except for the obvious $\eta \geq 0$ and no heuristic is performed.
- **h2:** We provide exactly one optimality cut to the initial MP in STEP 0 by solving the S-TSP with edge costs given by the design costs \bar{c}_e , $e \in E$. Such a solution also initializes the algorithm.
- **h3:** The initialization is as in h2. In addition, a modified two-opt local search procedure is applied to the initial solution to generate alternative solutions. The best γ solutions so obtained are used to generate optimality cuts. The best found solution initializes the algorithm.
- **h4:** The initialization is as for h3. The modified two-opt local search is applied every time an integer solution is found during the algorithm. The best γ solutions are used to generate optimality cuts and if improving solutions are found, the incumbent is updated.

Results for $\gamma = 5$ are summarized in Tables 3, 4, and 5 reporting separately the effect of instance size, demand type and value of α . In the three tables the first column reports the feature whose variation effects are examined. Columns two to five report information about the best upper bound obtained during the solution process. In particular, many experiments failed in obtaining a feasible solution when using h1, thus

the second column reports the percentage of such instances. We report in Column three, the average value of the upped bound obtained by h2. Columns four and five report the average relative improvements of h3 and h4 over h2, respectively, computed as $100(ub_2 - ub_i)/ub_2$, $i = 3, 4$. The remaining columns report the average computing times of h1, h2, h3, and h4, respectively.

Feat. dimension	Upper Bounds				Times (sec.)			
	h1 % no solution	h2	h3 %Impr.	h4 %Impr	h1	h2	h3	h4
40	30.0	573.8	0.2	0.4	10086.1	9370.0	9369.9	9372.9
50	45.0	645.1	0.3	0.4	13386.9	12475.3	12477.2	12470.3
60	61.7	692.4	0.3	0.4	15229.3	14040.3	14051.8	13960.8

Table 3: Comparison among heuristic strategies by instance dimension.

Feat. D. type	Upper Bounds				Times (sec.)			
	h1 % no solution	h2	h3 %Impr.	h4 %Impr	h1	h2	h3	h4
C	50.0	637.1	0.3	0.4	13421.9	12390.9	12396.4	12351.6
S	41.1	637.2	0.3	0.4	12379.6	11532.9	11536.2	11517.7

Table 4: Comparison among heuristic strategies by demand type

Feat. α	Upper Bounds				Times (sec.)			
	h1 % no solution	h2	h3 %Impr.	h4 %Impr.	h1	h2	h3	h4
0.25	3.3	883.1	0.0	0.0	4814.8	2544.7	2552.1	2453.8
0.5	38.3	636.9	0.0	0.0	15838.3	15289.5	15291.8	15294.5
0.75	95.0	391.3	0.8	1.1	18000.0	18000.0	18000.0	18000.0

Table 5: Comparison among heuristic strategies by α .

We observe that there is some advantage in adopting version h4. Indeed, in average it provides better upper bounds and computing times with respect to the other versions. This is due to the fact that, when h4 attains optimality, it is faster, and when it stops due to the time limit, it gives a better upper bound.

5.3.2 Optimality cuts and multicut generation.

As mentioned in the general description of the algorithm, optimality Benders cuts of the form (25) can be separated by solving a minimum cost flow problem for each commodity. Moreover, since at every master iteration all the subproblems share the same network, we can exploit sensitivity analysis to more efficiently obtain the solution of one subproblem given the solution of another.

It is moreover possible to exploit the decomposability of the Benders subproblem to produce a disaggregated version of the optimality cuts. Let η^{hk} be the contribution of the commodity (h, k) to the objective function, where $\sum_{(h,k) \in D} \eta^{hk} = \eta$. Substituting in (24), we obtain:

$$\min \alpha \sum_{(h,k) \in D} \eta^{hk} + (1 - \alpha) \sum_{e \in E} \bar{c}_e x_e, \quad (33)$$

and the disaggregated cuts are then

$$\eta^{hk} + \sum_{e \in E} \bar{\lambda}_e^{hk} x_e \geq \rho_h^{hk} - \rho_k^{hk} \quad \forall (h, k) \in D, \quad (34)$$

for every commodity (h, k) . Although disaggregated cuts can accelerate the convergence of the Benders algorithm, every time optimality cuts are separated, a potentially high number of cuts are added to the MP

and this can make MP considerably harder. For this reason, several authors proposed strategies to partially re-aggregate the optimality cuts (see, for example, Contreras et al. 2011).

We implemented and compared several versions of the optimality cuts generation. In all versions, the separation of the optimality cut is performed as described above, and the difference is in the different level of aggregation. At every separation of the optimality cut:

- **o1**: Only one optimality cut (25) is added.
- **o2**: For every origin destination (h, k) , a cut (33) is added.
- **o3**: We aggregate the optimality cuts of o2 by origin. More specifically, we consider η^h to be the contribution to the objective function of all commodities with origin in h , and we set $\sum_{h \in N} \eta^h = \eta$. Substituting in (24), we obtain:

$$\min \alpha \sum_{h \in N} \eta^h + (1 - \alpha) \sum_{e \in E} \bar{c}_e x_e, \tag{35}$$

and the corresponding partially disaggregated optimality cuts

$$\eta^h + \sum_{k \in N \setminus \{h\}} \sum_{e \in E} \bar{\lambda}_e^{hk} x_e \geq \sum_{k \in N \setminus \{h\}} \rho_h^{hk} - \rho_k^{hk} \quad \forall h \in N. \tag{36}$$

Results are summarized in Tables 6, 7, and 8. The first column of the three tables gives the instance feature considered by the table. The second column reports the root gap values which, as expected, are the same for all versions. In the next three columns, we report average computing times to reach optimality at the root node. The next three columns report average gaps obtained within five hours of computing time. The last three columns report the average computing times.

Feat.	root % gap	root time (sec.)			final % gap			total time (sec.)		
		o1, o2, o3	o1	o2	o3	o1	o2	o3	o1	o2
40	5.1	300.2	108.0	101.3	3.5	3.6	3.1	9370.0	9420.2	8600.8
50	6.1	1168.9	350.7	288.3	5.1	5.2	4.8	12475.3	12934.1	12343.2
60	6.1	2876.7	814.2	620.3	5.4	5.5	5.3	14040.3	14397.3	13777.4

Table 6: Comparison among Multicut strategies by instance dimension.

Feat.	root % gap	root time (sec.)			final % gap			total time (sec.)		
		D. type	o1, o2, o3	o1	o2	o3	o1	o2	o3	o1
C	5.7	2124.4	710.8	484.1	4.7	4.9	4.5	12390.9	13162.7	11864.3
S	5.9	772.8	137.9	189.2	4.6	4.6	4.3	11532.9	11338.4	11283.3

Table 7: Comparison among Multicut strategies by demand type.

Feat.	root % gap	root time (sec.)			final % gap			total time (sec.)		
		α	o1, o2, o3	o1	o2	o3	o1	o2	o3	o1
0.25	1.1	298.9	187.2	168.2	0.0	0.1	0.0	2544.7	3575.1	2205.1
0.5	4.1	898.7	352.9	294.4	2.5	2.7	2.1	15289.5	15162.5	14474.1
0.75	12.1	3148.2	732.8	547.2	11.5	11.5	11.1	18000.0	18000.0	18000.0

Table 8: Comparison among Multicut strategies by α .

We observe that o2 and especially o3 are much faster than o1 in obtaining root optimality. However, when looking at the overall solution process, this advantage is weaker. We observe that option o3 remains the best one, both for the gaps between lower and upper bounds, and the average computing times.

5.3.3 Separation strategies.

In our algorithms, we tried to keep the separation strategies as simple as possible. At every The algorithms check, at every separation round, for violated inequalities following the inequality check order. If violated inequalities of a given family are found, the search stops, all found inequalities are added to MP, and MP is reoptimized. When invoked again, the separation procedure will start from the beginning of the same order. The two sequences that we compared are:

1. **s1**: subtour elimination, heuristic blossom, exact blossom, combs from block, local cuts, optimality cuts.
2. **s2**: subtour elimination, optimality cuts, heuristic blossom, exact blossom, combs from block, local cuts.

The difference between s1 and s2 is the position of the optimality cuts: in s1 optimality cuts are only generated when all feasibility cuts are satisfied, while in s2, optimality cuts are separated when subtour inequalities are satisfied. In fact, when subtours elimination are satisfied, the network flow problem to separate optimality cuts is always feasible. The idea behind s2 is to avoid spending time trying to refine the fractional solution in a suboptimal region of the solution space.

Results are shown in Table 9, 10 and 11. As in the previous sections, the first column presents the instance feature. Columns two and three present the average gaps obtained at root by s1 and s2, respectively, and Columns four and five their computing times. Similarly, Columns six and seven report the average final gaps and, Columns 8 and 9 the average computing times of s1 and s2, respectively.

Feat.	root % gap		root time (sec.)		final % gap		final time (sec.)	
	s1	s2	s1	s2	s1	s2	s1	s2
dimension								
40	5.1	5.3	300.2	279.3	3.5	3.5	9370.0	9557.2
50	6.1	6.3	1168.9	1094.5	5.1	5.1	12475.3	12681.3
60	6.1	6.3	2876.7	2729.4	5.4	5.4	14040.3	14559.2

Table 9: Comparison between Separation Strategies by instance dimension.

Feat.	root % gap		root time (sec.)		final % gap		final time (sec.)	
	s1	s2	s1	s2	s1	s2	s1	s2
D. type								
C	5.7	5.9	2124.4	2028.4	4.7	4.7	12390.9	12696.3
S	5.9	6.0	772.8	707.1	4.6	4.6	11532.9	11835.6

Table 10: Comparison between Separation Strategies by instance type.

Feat.	root % gap		root time (sec.)		final % gap		final time (sec.)	
	s1	s2	s1	s2	s1	s2	s1	s2
α								
0.25	1.1	1.5	298.9	214.0	0.0	0.0	2544.7	3309.9
0.5	4.1	4.2	898.7	638.4	2.5	2.5	15289.5	15442.3
0.75	12.1	12.2	3148.2	3250.8	11.5	11.5	18000.0	18000.0

Table 11: Comparison between Separation Strategies by α .

Even though s1 and s2 should provide exactly the same results at the root node, we note that s1 yields better gaps than s2. The actual difference in lower bounds is due numerical instability of s2 that caused early exiting of the separation procedure. For this reason, we excluded s2 from further consideration, even if s1 is slower than s2 at the root.

5.3.4 Branching strategy.

Let (\bar{f}, \bar{x}) denote the optimal fractional solution of the current node. We considered the following two branching alternatives:

- **b1**: Branch on the variable the x_e closer to 0.5,
- **b2**: Branch on the variable x_e maximizing the following expression:

$$\min\{\bar{x}_e, 1 - \bar{x}_e\}(c_e + \sum_{(h,k) \in D} d_{hk} \bar{f}_{ij}^{hk}) \tag{37}$$

where $e = [i, j]$.

Strategy b1 is one of the most common ones, while strategy b2 tries to reroute the maximum possible amount of flow, with the hope of causing a stronger change in the subproblem’s objective value.

Results are reported in Tables 12, 13, and 14 where the first column identifies the instance feature, Columns two and three represent the average final gaps, Columns four and five the average computing times.

Tests mostly show that the considered branching strategies perform quite similarly in terms of gaps and computing times. However, b2 occasionally attains better gaps than b1 while showing comparable computing times. Consequently we chose b2 in our final experimentation.

Feat.	final % gap		final time (sec.)	
	b1	b2	b1	b2
dimension				
40	3.5	3.5	9370.0	9347.8
50	5.1	5.0	12475.3	12762.3
60	5.4	5.4	14040.3	13953.8

Table 12: Comparison between branching strategies: by dimension

Feat.	final % gap		final time (sec.)	
	b1	b2	b1	b2
D. type				
C	4.7	4.6	12390.9	12464.2
S	4.6	4.6	11532.9	11578.4

Table 13: Comparison between branching strategies: by demand density

Feat.	final % gap		final time (sec.)	
	b1	b2	b1	b2
α				
0.25	0.0	0.0	2544.7	2758.6
0.5	2.5	2.3	15289.5	15252.9
0.75	11.5	11.5	18051.5	18000.0

Table 14: Comparison between branching strategies: by α .

5.4 Final experiments

In accordance with the results of the experiments on the specific implementation details, we specialized the BBC3 algorithm by adopting heuristic h4, multicut generation o3, separation strategy s1, and branching rule b2. The resulting algorithm, called BBC3+, has been tested on the 48 benchmark instances of TSPGL2. Results are shown in Table 15. The first column indicates the name of the original TSPLIB instance where digits contained in the name give the number of nodes. The second column reports the demand type. Columns three and four report the percentage gap at the root node and the corresponding computing time. Similarly,

Name	D. type	root		final	
		GAP	Time	GAP	Time
att48	C	2.9	345.0	2.1	*
att48	S	4.8	133.9	3.3	*
bayg29	C	3.6	65.5	0.0	3242.0
bayg29	S	5.2	25.9	0.0	9433.5
berlin52	C	3.7	460.9	2.8	*
berlin52	S	1.6	164.3	0.5	*
burma14	C	1.5	2.2	0.0	4.7
burma14	S	1.8	1.6	0.0	4.9
dantzig42	C	3.2	125.2	0.0	14482.4
dantzig42	S	3.8	50.2	0.0	5807.9
eil51	C	6.2	498.0	5.6	*
eil51	S	3.7	151.1	2.1	*
eil76	C	5.1	3617.3	5.1	*
eil76	S	6.4	1212.1	6.3	*
eil101	C	7.8	12770.2	7.8	*
eil101	S	8.3	4250.7	8.2	*
fri26	C	1.9	24.4	0.0	103.0
fri26	S	1.9	13.4	0.0	83.6
gr17	C	1.7	3.0	0.0	16.5
gr17	S	3.3	1.6	0.0	60.3
gr21	C	0.9	7.0	0.0	8.9
gr21	S	2.0	4.1	0.0	13.8
gr24	C	1.3	19.8	0.0	55.6
gr24	S	3.3	7.7	0.0	385.8
gr48	C	6.1	437.0	5.0	*
gr48	S	5.4	139.9	3.3	*
gr96	C	6.4	8468.7	6.1	*
gr96	S	5.1	3331.6	5.0	*
kroA100	C	3.0	7173.2	2.8	*
kroA100	S	4.3	2434.0	3.9	*
kroB100	C	6.1	9642.2	6.0	*
kroB100	S	4.5	3569.5	4.3	*
kroC100	C	4.9	9061.9	4.7	*
kroC100	S	4.0	2573.2	3.7	*
kroD100	C	4.4	7486.5	4.2	*
kroD100	S	4.5	3307.8	4.4	*
pr76	C	7.1	1973.5	7.0	*
pr76	S	4.2	704.7	3.8	*
rat99	C	3.9	4708.9	3.7	*
rat99	S	8.0	1821.4	7.9	*
rd100	C	5.2	9468.4	4.9	*
rd100	S	4.8	3088.0	4.6	*
st70	C	4.7	2149.1	4.7	*
st70	S	4.1	748.1	3.4	*
ulysses16	C	2.4	5.3	0.0	91.9
ulysses16	S	1.6	3.3	0.0	24.7
ulysses22	C	1.2	15.3	0.0	103.9
ulysses22	S	3.6	4.7	0.0	1936.6

Table 15: Final experimentation on TSPLIB-derived instances, $\alpha = 0.5$

Columns five and six report the final gap and the overall computing time. A sign “*” in the last column indicates that the algorithm did not prove optimality within 5 hours of computing time.

We observe that BBC3+ was able to optimally solve 18 of the 48 instances. The largest solved instance is *dantzig42* with complete demand matrix. The hardest considered instances turned out to be *eil101* with sparse demand matrix. BBC3+ was not able to prove optimality for such an instance, but the lower bound guarantees that the best feasible solution value is within 8.2 percent of the optimal solution. The average best feasible solution value computed on the set of problems not solved to optimality is within 4.6 percent of the optimal solution in the worst case.

As mentioned in Section 5.1, the considered instance dimension range in TSPGL2 is representative of typical situations in our reference application. We observe that from such a viewpoint, BBC3+ can be also used as a very good heuristic, given the quite low optimality gaps provided on instances not solved to optimality.

6 Conclusions

We presented a new problem in the class of TSP conjugating aspects typical in the TSP field and aspects coming from network design and routing. The conjunction of those aspects makes the problem extremely challenging even for relatively small instances. The problem has different applications, one of which is within the design of semi-flexible transit services. We proposed a multicommodity flow formulation and derived a particular branch-and-cut algorithm based on Benders reformulation. This approach takes advantage of important properties that relate our problem feasible region with the TSP polyhedron. An extensive computational experimentation compares some variants of the proposed algorithm and, for small instances only, a commercial solver. We also introduce a set of benchmark instances derived from TSPLIB, that in some cases constitute an interesting challenge for future research. Among the future aspects to be explored there is the study of specific relaxation strategies with the intent of tightening the current lower bound that are not too tight especially when the routing component of the problem is more important than the design component.

Acknowledgement

While working on this paper, F. Errico was postdoctoral researcher at CIRRELT and École des sciences de la gestion, Université du Québec à Montréal.

T.G. Crainic was the NSERC Industrial Research Chair in Logistics Management, ESG UQAM, and Adjunct Professor with the Department of Computer Science and Operations Research, Université de Montréal, and the Department of Economics and Business Administration, Molde University College, Norway, and F. Errico was postdoctoral researcher with the Chair.

Funding for this project was provided by the Natural Sciences and Engineering Council of Canada (NSERC), through its Industrial Research Chair and Discovery Grant programs, by our partners CN, Rona, Alimentation Couche-Tard, the Ministry of Transportation of Québec, and by the Fonds de recherche du Québec - Nature et technologies (FRQNT) through its Team Research Project program.

We also gratefully acknowledge the support of the Canada Foundation for Innovation, the Québec Ministry of Education, UQAM and SUN/Azuris for our computing infrastructure, as well as that of Fonds de recherche du Québec for the CIRRELT infrastructure grants.

References

- Applegate, D., R. Bixby, V. Chvátal, and W. Cook (2003). Concorde: Traveling Salesman Problem Solver. <http://www.tsp.gatech.edu/concorde.html>.
- Applegate, D., R. Bixby, V. Chvátal, and W. Cook (2007). *The Traveling Salesman Problem: A Computational Study*. Princeton University Press.
- Atamtürk, A. (2002). On capacitated network design cut-set polyhedra. *Mathematical Programming* 92, 425–437.
- Barahona, F. (1996). Network design using cut inequalities. *SIAM Journal on Optimization* 6(3), 823–837.
- Benders, J. (1962). Partitioning procedures for solving mixed-variables programming problems. *Numerische Mathematik* 4.
- Bienstock, D. and O. Günlük (1994). Capacitated Network Design – Polyhedral Structure and Computation. *INFORMS Journal on Computing* 8, 243–259.
- Contreras, I., J.-F. Cordeau, and G. Laporte (2011). Benders decomposition for large-scale uncapacitated hub location. *Operations Research* 59(6), 1477–1490.
- Cordeau, J.-F., F. Soumis, and J. Desrosiers (2000). A Benders Decomposition Approach for the Locomotive and Car Assignment Problem. *Transportation Science* 34(2), 133–149.
- Cordeau, J.-F., F. Soumis, and J. Desrosiers (2001). Simultaneous Assignment of Locomotives and Cars to Passenger Trains. *Operations Research* 49(4), 531–548.
- Costa, A. (2005). A survey on benders decomposition applied to fixed-charge network design problems. *Computers & Operations Research* 32(6), 1429–1450.
- Crainic, T. G., A. Frangioni, and B. Gendron (2001). Bundle-based relaxation methods for multicommodity capacitated fixed charge network design. *Discrete Applied Mathematics* 112(1-3), 73 – 99.
- Desaulniers, G. and M. D. Hickman (2007). Public Transit. In Barnhart, C. and Laporte, G. (Eds.), *Transportation*, Volume 14 of *Handbooks in Operations Research and Management Science*, pp. 69–127. North-Holland.
- Errico, F. (2008). *The design of flexible transit systems: models and solution methods*. Ph. D. thesis, Politecnico di Milano, Milano, Italy.
- Errico, F., T. G. Crainic, F. Malucelli, and M. Nonato (2011a). A unifying framework and review of Semi-Flexible Transit Systems. CIRRELT-2011-64, Centre interuniversitaire de recherche sur les réseaux d'entreprise, la logistique et le transport, Université de Montréal.
- Errico, F., T.G. Crainic, F. Malucelli, and M. Nonato (2011b). The Design Problem for Single-line Demand Adaptive transit Systems. CIRRELT-2011-65, Centre interuniversitaire de recherche sur les réseaux d'entreprise, la logistique et le transport, Université de Montréal.
- Fischetti, M., G. Laporte, and S. Martello (1993). The Delivery Man Problem and Cumulative Matroids. *Operations Research* 41(6), 1055–1064.
- Günlük, O. (1999). A branch-and-cut algorithm for capacitated network design problems. *Mathematical Programming* 86, 17–39.

- Holmberg, K. and J. Hellstrand (1998). Solving the Uncapacitated Network Design Problem by a Lagrangean Heuristic and Branch-and-Bound. *Operations Research* 46(2), 247–259.
- Holmberg, K. and D. Yuan (2000). A Lagrangian Heuristic Based Branch-and-Bound Approach for the Capacitated Network Design Problem. *Operations Research* 48(3), 461–481.
- Kliwer, G. and L. Timajev (2005). Relax-and-Cut for Capacitated Network Design. In G. Brodal and S. Leonardi (Eds.), *Algorithms - ESA 2005*, Volume 3669 of *Lecture Notes in Computer Science*, pp. 47–58. Springer Berlin / Heidelberg.
- Leung, J. M. Y. and T. L. Magnanti (1989). Valid inequalities and facets of the capacitated plant location problem. *Mathematical Programming* 44, 271–291.
- Ljubić, I., P. Putz, and J.-J. Salazar-González (2012). Exact approaches to the single-source network loading problem. *Networks* 59(1), 89–106.
- Lucena, A. (1990). Time-dependent Traveling Salesman Problem - The Deliveryman case. *Networks* 20(6), 753–763.
- Magnanti, T. and R. Wong (1981). Accelerating benders decomposition: Algorithmic enhancement and model selection criteria. *Operations Research* 29(3).
- Magnanti, T. L., P. Mireault, and R. T. Wong (1986). Tailoring Benders decomposition for uncapacitated network design. In *Netflow at Pisa*, Volume 26 of *Mathematical Programming Studies*, pp. 112–154. Springer Berlin Heidelberg.
- McDaniel, D. and M. Devine (1977). A Modified Benders' Partition Algorithm for Mixed Integer Programming. *Management Science* 24(3), 312–319.
- Méndez-Díaz, I., P. Zabala, and A. Lucena (2008). A new formulation for the Traveling Deliveryman Problem. *Discrete Appl. Math.* 156(17), 3223–3237.
- Naddef, D. and S. Thienel (2002). Efficient separation routines for the symmetric traveling salesman problem I: General tools and comb separation. *Mathematical Programming* 92, 272–255.
- Naoum-Sawaya, J. and S. Elhedhli (2012). An interior-point benders based branch-and-cut algorithm for mixed integer programs. *Annals of Operations Research*, 1–23. 10.1007/s10479-010-0806-y.
- Ortega, F. and W. L. A. (2003). A branch-and-cut algorithm for the single-commodity, uncapacitated, fixed-charge network flow problem. *Networks* 41(3), 143–158.
- Padberg, M. and M. R. Rao (1982). Odd minimum cut-set and b -matchings. *Mathematics of Operations Research* 7, 62–80.
- Padberg, M. and G. Rinaldi (1990). An Efficient Algorithm for the Minimum Capacity Cut Problem. *Mathematical Programming* 47, 19–36.
- Padberg, M. and T. Sung (1991). An analytical comparison of different formulations of the travelling salesman problem. *Mathematical Programming* 52(1), 315–357.
- Potts, J. F., M. A. Marshall, E. C. Crockett, and J. Washington (2010). A Guide for Planning and Operating Flexible Public Transportation Services. Volume 140 of *TCRP Synthesis Report*. Transportation Research Board.

- Raack, C., A. M. Koster, S. Orlowski, and R. Wessälly (2011). On cut-based inequalities for capacitated network design polyhedra. *Networks* 57(2), 141–156.
- Rei, W. and Cordeau, J.-F., M. Gendreau, and P. Soriano (2009). Accelerating Benders Decomposition by Local Branching. *INFORMS Journal on Computing* 21(2), 333–345.
- Rodríguez-Martín, I. and J.-J. Salazar-González (2008). Solving a capacitated hub location problem. *European Journal of Operational Research* 184(2), 468 – 479.
- Sarubbi, J. and H. Luna (2007). The Multicommodity Traveling Salesman Problem. In *Proceedings of International Network Optimization Conference*. Spa, Belgium.
- Sarubbi, J., G. Miranda, H. Luna, and G. Mateus (2008). A Cut-and-Branch algorithm for the Multicommodity Traveling Salesman Problem. In *IEEE international conference on Service Operations and Logistics, and Informatics*, pp. 1806–1811. IEEE/SOLI, Beijing.
- Sellmann, M., G. Kliewer, and A. Koberstein (2002). Lagrangian Cardinality Cuts and Variable Fixing for Capacitated Network Design. In *Proceedings of the 10th Annual European Symposium on Algorithms, ESA '02*, London, UK, UK, pp. 845–858. Springer-Verlag.
- Sridhar, V. and J. S. Park (2000). Benders-and-cut algorithm for fixed-charge capacitated network design problem. *European Journal of Operational Research* 125(3), 622 – 632.
- TSPGL2 (2012). The Symmetric Traveling Salesman Problem with Generalized Latency: the TSPGL2 instance set. <http://w1.cirrelt.ca/~errico/TSPGL2.zip>.