

Computing the Envelope for Stepwise-Constant Resource Allocations

Nicola Muscettola
 NASA Ames Research Center
 Moffett Field, California 94035-1000
mus@email.arc.nasa.gov

Abstract. Computing tight resource-level bounds is a fundamental problem in the construction of flexible plans with resource utilization. In this paper we describe an efficient algorithm that builds a resource envelope, the tightest possible such bound. The algorithm is based on transforming the temporal network of resource consuming and producing events into a flow network with nodes equal to the events and edges equal to the necessary predecessor links between events. A staged maximum flow problem on the network is then used to compute the time of occurrence and the height of each step of the resource envelope profile. Each stage has the same computational complexity of solving a maximum flow problem on the entire flow network. This makes this method computationally feasible and promising for use in the inner loop of flexible-time scheduling algorithms.

1 Resource Envelopes

Retaining temporal flexibility in activity plans is important for dealing with execution uncertainty. For example, flexible plans allow explicit reasoning about the temporal uncontrollability of exogenous events [11] and the seamless incorporation of execution countermeasures. Fixed-time schedules (i.e., the assignment of a precise start and end time to all activities) are brittle and it is typically very difficult to exactly follow them during execution. For an example of the effect of fixed-time schedules in an intelligent execution situation, consider the “Skylab strike” [6], when during the Skylab 4 mission, astronauts went on a sit-down strike after 45 days of trying to catch up with the exact demands of a fast paced schedule with no flexibility for them to adjust to the space environment.

A major obstacle to using flexible schedules, however, remains the difficulty of computing the amount of resources needed across all of their possible executions. This problem is particularly difficult for multiple-capacity resources (such as electrical power) that can be both consumed and produced in any amount by concurrent activities. Techniques have been developed [5] [10] for giving conservative estimates of the resource levels needed by a flexible schedule, yielding both an upper bound and a lower bound profile on the resource level over time. In the context of a systematic search method to build flexible plans, resource-level bounds can be used at each search step as follows: a) as a backtracking test, i.e., to determine when the lower/upper bound interval is outside of the range of allowed resource levels at some time and therefore no fixed-time instantiations of the plan is resource-feasible; and b) as a search termination test, i.e., to determine when the lower/upper bound range is inside the range of allowed resource levels at all times and therefore all fixed-time instantiations of the plan are resource-feasible.

Bound tightness is extremely important computationally since a tight bound can save a potentially exponential amount of search (through early backtracking and solution detection) when compared to a looser bound. In this paper, we discuss how to compute the *resource-level envelope*, i.e., the measure of maximum and minimum resource consumption at any time for all fixed-time schedules in the flexible plan. At each time the envelope guarantees that there are

two fixed-time instantiations, one producing the minimum level and the other the maximum. Therefore, the resource-level envelope is the tightest possible resource-level bound for a flexible plan since any tighter bound would exclude the contribution of at least one fixed-time schedule. If the resource-level envelope can be computed efficiently, it could substitute looser bounds that are currently used in the inner core of constraint-posting scheduling algorithms (Laborie 2001) with the potential for great improvements in performance.

To appreciate the difficulty of computing the resource level, we can compare the cases of a fully flexible plan with that of a plan with a single fixed-time instantiation. In the fixed-time case, the envelope degenerates to the resource profile that is used in the inner loop of traditional scheduling algorithms [2][15]. Computing from scratch a resource profile is cheap. It is easy to see that its worst case time complexity is $O(N \lg N)$ where N is the number of activities in the flexible plan. Consider now a fully flexible plan and assume that we naively wanted to compute the resource-level envelope for a flexible plan by simply enumerating all schedules and taking the maximum/minimum of all resource levels at all times. Since a flexible activity plan has a number of possible instantiated fixed-time schedules that is exponential in the number of events, such naïve method is clearly impractical in most cases.

Note that “resource-level envelope calculation” is *not* equivalent to “scheduling with multiple capacity”, an NP-hard problem. For example, note that a reduction of the envelope calculation to scheduling is not straightforward. Consider the interval between minimum and maximum envelope at a given time. It is true that a solution to the scheduling problem can be obtained in polynomial time if the envelope interval is always completely contained within resource availability (in which case all fixed-time instantiations are legal schedules). However, if at least two times the envelope interval is partly outside the availability bound, one can generate examples in which the envelope cannot tell us whether there is a fixed-time schedule that is within availability at all times or whether for all schedules the resource level is outside availability at some time. Discriminating between these two cases still requires search.

This paper presents a polynomial algorithm for the computation of resource-level envelopes based on a novel combination of the theory of shortest-paths in the temporal constraint network for the flexible plan, and the theory of maximum flows for a flow network derived from the temporal and resource constraints. We develop the theory, show that the algorithm is correct, and that its asymptotic complexity is $O(N O(\text{maxflow}(N)))$, where N is the number of start/end times of the activities in the plan, which is at most twice the number of activities, and $O(\text{maxflow}(N))$ is the complexity of a maximum flow algorithm applied to an auxiliary flow network with N nodes. We believe that this method will be efficient in practice, since experimental analysis [1] show the practical cost of maxflow to be as good as $O(N^{1.5})$. However this paper is a theoretical contribution and a definitive answer to its practical complexity will require further experimental work.

In the rest of the paper we introduce some notation and describe the formal model of activity networks with resource consumption. Then we review the literature on resource contention measures and show an example in which the current state of the art is inadequate. The discussion of our algorithm follows. Some informal examples to establish an intuitive understanding of our method are first given. Then we establish the connection between maximum flow problems and finding sets of activities that have the optimal contribution to the resource level. These sets are then shown to compute an envelope. Finally, we describe a simple envelope algorithm and its complexity, and conclude discussing future work.

2 Resource Envelopes

Figure 1 shows an activity network with resource allocations. Our notation is equivalent to that of previous work on flexible scheduling with multiple-capacity resources [5][10]. The network has two time variables per activity, a start event and an end event (e.g., e_{1s} and e_{1e} for activity A_1), a non-negative flexible activity duration link (e.g., $[2, 5]$ for activity A_1), and flexible

separation links between events (e.g., $[0, 4]$ from e_{3e} to e_{4e}). A time origin, T_s , corresponds to time 0 and supports separation links to other events. Without loss of generality we assume that all events occur after T_s and before an event T_e rigidly connected to T_s . The interval $T_s T_e$ is the *time horizon T* of the network.

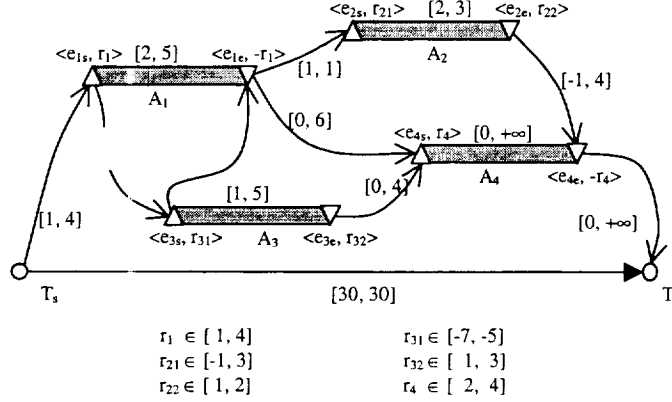


Figure 1: An activity network with resource allocations

Time origin, events and links constitute a Simple Temporal Network (STN) [8]. Unlike regular STNs, however, each event has an associated *allocation variable* with real domain (e.g., r_{31} for event e_{3e}) representing the amount of resource allocated when the event occurs. We call this augmented network **R** a piecewise-constant Resource allocation STN (cR-STN). In the following we will assume that all allocations refer to a single, multi-capacity resource. The extension of the results to the case of multiple resources is straightforward. An event e^- with negative allocation is a *consumer*, while an e^+ with positive allocation is a *producer*.

This formalization covers many of the usual models for resource allocations. For example, note that an event can be either a consumer or a producer in different instantiations of the allocation variables (e.g., event e_{2s} for which the bound for r_{21} is $[-1, 3]$). This allows reasoning about dual-use activities (e.g., the activities of starting a car and running it make use of the alternator as either a power consumer or producer). Moreover, some events can have resource allocations that are opposite to each other (e.g., e_{1e} vs. e_{1s}). This allows modeling allocations that last only during an activity's occurrence, such as power consumption. Note, however, that this model does not cover continuous accumulation such as change of energy stored in a battery over time. A conservative approximation can be achieved by accounting for the entire resource usage at the activity start or end.

We will always assume that the cR-STN is temporally consistent. From the STN theory [8], this means that the shortest-path problem associated with **R** has a solution. Given two events e_1 and e_2 we denote with $|e_1 e_2|$ the shortest-path from e_1 to e_2 . We will call a full instantiation of the time variables in **R** a *schedule* $s(\cdot)$ where $s(e)$ is the time of occurrence of event e according to schedule s . The set **S** contains all possible consistent schedules for **R**. Each event e has a time bound $[et(e), lt(e)]$, with $et(e)$ and $lt(e)$ respectively the earliest and latest time for e . The time bound represents the range of possible time values $s(e)$ for all $s \in \mathbf{S}$. From the STN theory, we know that $et(e) = -|e T_s|$ and $lt(e) = |T_e e|$. Finally, given three events, e_1 , e_2 and e_3 , the triangular inequality among shortest paths $|e_1 e_3| \leq |e_1 e_2| + |e_2 e_3|$ holds.

A fundamental data structure used in the rest of the paper is the *anti-precedence graph*, **Aprec(R)**, for a cR-STN **R**. The anti-precedence graph is similar to the precedence graph of [10] with the following differences: 1) the links are in reverse order; 2) it does not distinguish between the set of strict precedence edges, $E_<$, and the set of precedence edges that allow time equality, $E_≤$; and 3) several possible kinds of precedence graphs are allowable for a network **R**. Formally, **Aprec(R)** is a graph with the same events as **R** and such that for any two events

e_1 and e_2 with $|e_1 e_2| \leq 0$ there is a path from e_1 to e_2 in $\text{Aprec}(\mathbf{R})$. Alternatively, we can say that an event e_1 precedes another e_2 in the anti-precedence graph if e_1 cannot be executed before e_2 . A way to build an anti-precedence graph is to run an all-pairs shortest-path algorithm on \mathbf{R} and retain only the edges with non-positive shortest distance. Smaller graphs can also be obtained by eliminating dominated edges. The choice of precedence graph type may affect performance but not the correctness of the algorithm described here.

The cost of computing $\text{Aprec}(\mathbf{R})$ is therefore bound by the cost of computing the all-pairs shortest path graph for \mathbf{R} , i.e., it is $O(NE + N^2 \lg V)$ where N is the number of events (at most twice the number of activities in the plan) and E is the number of temporal distance constraints in the original cR-STN [7].

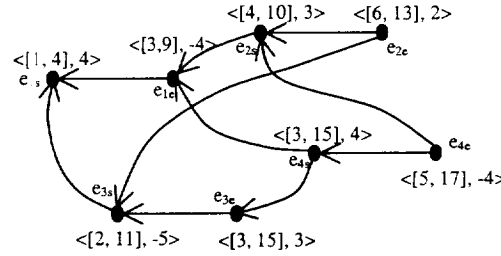


Figure 2: Anti-precedence graph with time/resource usage

Figure 2 depicts one of the anti-precedence graphs of the network in Figure 1 with the time bound and the maximum allowed resource allocation labeling each event.

3 Resource Contention Measures

Safe execution of a flexible activity networks needs to avoid resource contention, i.e., the possibility that for some consistent time assignment to the events there is at least one time at which the total amount of resource allocated is outside the availability bounds. There are essentially two methods for estimating resource contention: heuristic and exact. Most of the heuristic techniques [14] [12] [3] measure the probability of an activity requesting a resource at a certain time. This probability is estimated either analytically on a relaxed constraint network or stochastically by sampling time assignments on the full constraint network. The occurrence probabilities are then combined in an aggregate demand on resources over time, the contention measure. Probabilistic contention can give a measure of likelihood of a conflict. However, it is not a safe measure, i.e., the lack of detected conflicts does not exclude that a variable instantiation for the cR-STN could cause an inconsistent resource allocation.

Exact methods avoid this problem. They compute sufficient conditions for the lack of contention. [10] has a good survey of such methods. Current exact methods operate on relaxations of the constraint network. For example, edge-finding techniques [13] analyze how an activity can be scheduled relative to a subset of activities, comparing the sum of all durations with a time interval derived from the time bounds of all the activities under consideration. Relying solely on time bounds ignores much of the inter-activity constraints and is effective only when the time bounds are relatively tight. Therefore algorithms using these contention measures tend to eliminate much of the flexibility in the activity network. Some recent work [5] [10] goes further in exploiting inter-activity constraints. For example, [10] proposes a *balance constraint* that is based on conservative upper and lower bounds on the resource level immediately before and after each event e . These bounds precisely estimate the contribution of events that must precede e and overestimate the contribution of events that may or may not precede e . The over-estimate assumes that only the events with the worst contribution (producers for upper bounds and consumers for lower bounds) happen before e . The balance constraint appears to

work well in a flexible scheduling algorithm [10] but the bounds on which it is based may be very loose for networks with significant amounts of parallelism.

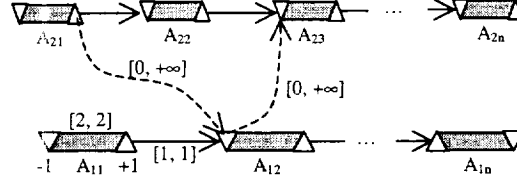


Figure 3: Over-constraining a flexible activity network

For example, consider the activity graph in Figure 3, consisting of two rigid chains of n activities with the same fixed duration and the same fixed activity separation. Assume that the horizon T is wide enough to allow any feasible ordering among them. Each activity consumes one unit during its occurrence and the resource has two available units of capacity over time. It is clear that the maximum resource level requested by the flexible plan at any time cannot exceed two, and therefore all schedules obtained by merging any schedule of the two chains are feasible. A scheduling algorithm using a resource-level envelope will therefore recognize that the initial problem is already a feasible flexible plan and will terminate immediately. In contrast, a scheduler using the balance constraint will always detect an over-allocation until it somehow constrains the network (e.g., by systematic or local search) with constraints that are at least as tight as the two following cases: a) the start activity n of one chain occurs no later than the start of the second activity of the other; or b) more than two activities overlap and there is an activity k on one chain that must start between the end of activity i and the start of activity $i+2$ on the other chain. The dashed arrows in Figure 3 represent the constraints posted in case b. The balance constraint cannot correctly handle this situation because it cannot account for the constraint structure of all possible parallel chains *simultaneously*. The rest of this paper shows that the full constraint structure can be efficiently exploited in calculating the resource-level envelope.

4 Resource Envelopes

As discussed in the introduction, we are seeking the maximum and minimum possible resource production (consumption) among all possible schedules of R . Note that the maximum (minimum) overall resource level induced by R for any possible schedule can always be obtained by assigning each allocation variable to its maximum (minimum) possible value. For any specific value assignment to the allocation variables, each event has a constant weight: positive, $c(e^+)$, for a producer and negative, $-c(e^-)$, for a consumer. Given a schedule $s \in S$ and a time $t \in T$, $E_s(t)$ is the set of events e such that $s(e) \leq t$. For any subset A of the set of events in R , $E(R)$, we define the *resource-level increment* as $\Delta(A) = 0$ if $A = \emptyset$, and $\Delta(A) = \sum_{e^+, e^- \in A} c(e^+) - c(e^-)$ if $A \neq \emptyset$. The following functions of time rigorously define the resource-level envelope:

- *Resource level* due to schedule s : $L_s(t) = \Delta(E_s(t))$.
- *Maximum resource envelope*: $L_{\max}(t) = \max_{s \in S} (L_s(t))$.
- *Minimum resource envelope*: $L_{\min}(t) = \min_{s \in S} (L_s(t))$.

The *resource level envelope* that we seek is the interval bound function of time $[L_{\min}(t), L_{\max}(t)]$. Since the methods to compute L_{\min} and L_{\max} can be obtained from each other with obvious term substitutions, we only develop the algorithm for L_{\max} .

Before we formally discuss our algorithm, we want to introduce some examples to give an intuitive feel of the foundations of the method. First consider an activity network consisting of

a single activity that produces or consumes resource capacity during its occurrence (Figure 4(a)). We could build L_{\max} by asking at each time $t \in T$ whether A_1 can happen before, after or can overlap t . If the activity starts with a resource production (Figure 4(b)), then the resource level will be maximum at t if A_1 starts, contains or ends at t . This is always possible between $et(e_{1s})$ and $lt(e_{1e})$. Within this interval $L_{\max}(t) = 1$, while outside $L_{\max}(t) = 0$. Conversely, if A_1 starts with a consumer (Figure 4(c)), then the maximum resource level can be zero at time t only if A_1 can start after t or can end before t for some schedule. This is possible only before $lt(e_{1s})$ and after $et(e_{1e})$. Therefore, $L_{\max}(t) = -1$ between $lt(e_{1s})$ and $et(e_{1e})$, and $L_{\max}(t) = 0$ everywhere else. This example suggests a strategy for computing L_{\max} that looks at each event and considers the incremental contribution of the event's weight to the maximum resource envelope at the earliest time for producers or at the latest time for consumers.

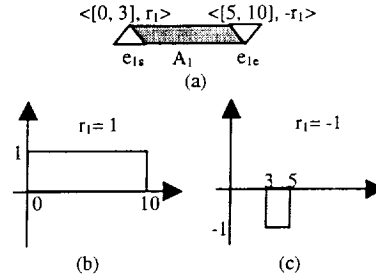


Figure 4: Maximum resource-level envelope for a single activity

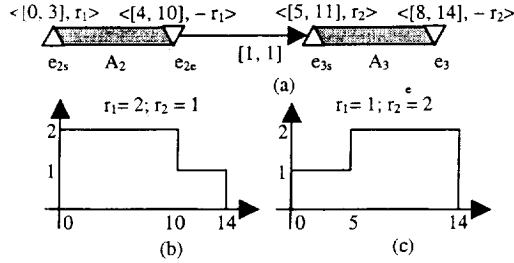


Figure 5: Maximum level envelope for two chained activities

For a complex network, however, this simple strategy is insufficient. Consider a rigidly linked pair of activities with a reusable resource allocation (Figure 5(a)). In this case, the time of occurrence of e_{2e} and e_{3s} are bound together. Looking at the contribution to the envelope of each event in isolation, we would want to add the contribution of e_{2e} as late as possible since it is a consumer, and the contribution of e_{3s} as early as possible, since it is a producer. The decision of which time to choose depends on the total contribution of *both* events. The total contribution will be added at $lt(e_{2e})$ if the total contribution is a consumption (Figure 5(b)) or at $et(e_{3s})$ if the total contribution is a production (Figure 5(c)). Note that in both cases e_{2s} and e_{3e} are *pending* at the selected time, i.e., their contribution has not been added yet to the envelope but they both could occur at the selected time. This suggests revising the strategy for computing L_{\max} as follows: at the earliest or latest time of each event, select the set of pending events whose resource-level increment is maximum, eliminate these events from the pending events and declare the increment for L_{\max} at that time to be their resource-level increment.

Now consider the network in Figure 1 and the event time bounds, maximum resource allocation and precedence graph in Figure 2. Assume that we want to compute $L_{\max}(3)$. The set of events that may be scheduled before, at or after time 3 is $\{e_{1s}, e_{1e}, e_{3s}, e_{3e}, e_{4s}\}$. However, of these only $\{e_{1e}, e_{3s}, e_{3e}, e_{4s}\}$ are pending since the contribution to L_{\max} of e_{1s} occurs at its earliest time 1. The subset of pending events that we need to consider at time 3 to compute the

increment on L_{\max} are all those that are forced to occur at or before 3 assuming that some pending event occurs at 3. These subsets are $\{e_{1e}\}$, $\{e_{3s}\}$, $\{e_{3s}, e_{3e}\}$ and $\{e_{1e}, e_{3s}, e_{3e}, e_{4s}\}$. Unfortunately, each of these subsets has a negative weight and therefore none of them is selected at time 3 since, from the definition of resource-level increment, the empty event set has the maximum (zero) increment. At time 4 the set of pending events is augmented with e_{2s} and a new subset of pending events with positive weight, $\{e_{1e}, e_{2s}, e_{3s}, e_{3e}, e_{4s}\}$, is possible. This generates the increment that, added to $L_{\max}(3)$, gives us $L_{\max}(4)$.

The selection of the pending events subset of maximum resource-level increment is the key source of complexity in calculating L_{\max} . An exhaustive enumeration of all subsets is intractable in the general case. Fortunately, it turns out that this selection problem is equivalent to a maximum flow problem for an appropriate flow network derived from $\text{Aprec}(\mathbf{R})$. We discuss this rigorously in the rest of the paper.

5 Calculating Maximum Resource-Level Increments

Consider an interval $H \subseteq T$. We can partition all events in \mathbf{R} into three sets depending on their relative position with respect to H : 1) the *closed events* C_H with all events that must occur strictly before or at the start of H , i.e., such that $lt(e) \leq \text{start}(H)$; 2) the *pending events* R_H with all events that can occur within or at the end of interval H , i.e., such that $lt(e) > \text{start}(H)$ and $et(e) \leq \text{end}(H)$; and 3) the *open events* O_H with all events that must occur strictly after H , i.e., such that $et(e) > \text{end}(H)$.

The set R_H could contain events that can be scheduled both inside and outside H . If $H=T$, then $C_T = \emptyset$, $R_T = E(\mathbf{R})$ and $O_T = \emptyset$. If H is a single instant of time, i.e., $H=[t, t]$, we will use the simplifying notation $C_t = C_{[t, t]}$, $R_t = R_{[t, t]}$ and $O_t = O_{[t, t]}$.

Assume that we want to compute the resource-level increment for a schedule s at a time $t \in H$. This will always include the contribution of all events in C_H and none of those in O_H irrespective of s and t . With respect to the events in R_H , if an event is scheduled to occur at or before t , then all of its predecessors (according to $\text{Aprec}(\mathbf{R})$) will also have to occur at or before t . In other words, it is possible to find a set of events $X \in R_H$ such that the events $e_p \in R_H$ that are scheduled no later than t in s are those in X or those $e_p \in R_H$ such that $|e_x, e_p| \leq 0$ for some $e_x \in X$. We call this the *predecessor set* P_X of X . Therefore, the resource level at time t for a given schedule s is the sum of the weights of events in C_H and in $P_X \subseteq R_H$. Since we are trying to maximize the resource level, we will look for the sets P_X with the maximum total weight. An important property that we will exploit later is that given two predecessor sets P_X and P_Y , $P_X \cap P_Y$ and $P_X \cup P_Y$ are also predecessor sets.

5.1 Resource-Level Increments and Maximum Flow

We know that to compute L_{\max} we want to look for sets P_X with maximum weight. We find these by computing a maximum flow for an auxiliary flow network built from R_H and $\text{Aprec}(\mathbf{R})$. For a complete discussion of maximum flow problems see [7]. Here we only highlight some concepts and relations that we will use.

First let us define the auxiliary flow problem that we will use to compute L_{\max} .

Resource Increment Flow Problem: Given a set of pending events R_H for a cr -STN \mathbf{R} , we define the resource increment flow problem $F(R_H)$ with source σ and sink τ as follows:

1. For each event $e \in R_H$ there is a corresponding node $e \in F(R_H)$.
2. For each event $e^+ \in R_H$, there is an edge $\sigma \rightarrow e^+$ with capacity $c(e^+)$.
3. For each event $e^- \in R_H$, there is an edge $e^- \rightarrow \tau$ with capacity $c(e^-)$, i.e., the opposite of e^- 's weight in \mathbf{R} .
4. For each pair of e_1 and e_2 with an edge $e_1 \rightarrow e_2$ in the anti-precedence graph $\text{Aprec}(\mathbf{R})$, there is a corresponding link $e_1 \rightarrow e_2$ in $F(R_H)$ with capacity $+\infty$.

Figure 6 shows the auxiliary flow problem for the anti-precedence graph in Figure 2, with every edge labeled with its capacity.

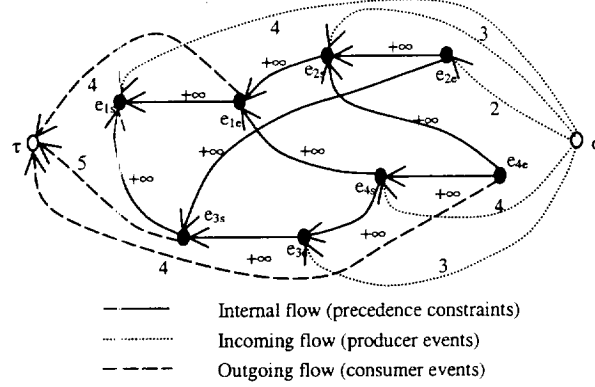


Figure 6: Resource increment flow problem

Consistent with the theory of maximum flows, we will indicate $f(e_1, e_2)$ as the flow associated to a link $e_1 \rightarrow e_2$ in $F(R_H)$. The flow function has, by definition, several properties. It is skew-symmetric, i.e., $f(e_2, e_1) = -f(e_1, e_2)$. Each flow has to be not greater than the capacity of the link to which it is associated. For example, referring to the flow network in Figure 6, $f(\sigma, e_{2e}) \leq 2$, while $f(e_{2e}, \sigma) \leq 0$ since there is no edge from e_{2e} and σ , a situation equivalent to the capacity of the edge $e_2 \rightarrow \sigma$ being zero. We also use the implicit summation notation $f(A, B) = \sum_{a \in A} \sum_{b \in B} f(a, b)$, where B and A are disjoint event sets in $F(R_H)$, to indicate the flow $f(A, B) = \sum_{a \in A} \sum_{b \in B} f(a, b)$. Consider now any subset of events $A \subseteq R_H$ and let us call A^c the set of events $A^c = R_H - A$. From the final property defining a flow function, flow conservation, we can obtain the following: $f(\{\sigma\}, A) = f(A, \{\tau\}) + f(A, A^c)$. The total network flow is defined as $f(\{\sigma\}, R_H) = f(R_H, \{\tau\})$. The maximum flow of a network is a flow function f_{\max} such that the total network flow is maximum.

A fundamental concept in the theory of flows is the *residual network*. This is a graph with an edge for each pair of nodes in $F(R_H)$ with positive *residual capacity*, i.e., the difference between edge capacity and flow. Each edge in the residual network has capacity equal to the residual capacity. For example, considering the network in Figure 6, assume that $f(e_{1e}, \tau) = 3$ and $f(\sigma, e_{2e}) = 2$. The residual network for that flow will have the following edges: $e_{1e} \rightarrow \tau$ with capacity 1, $\tau \rightarrow e_{1e}$ with capacity 3, and $e_{2e} \rightarrow \sigma$ with capacity 2. Also note that any residual network for any flow of $F(R_H)$ will always have an edge of infinite capacity for each edge in the precedence graph $\text{Aprec}(R)$. An *augmenting path* is a path connecting σ to τ in the residual network. The existence of an augmenting path indicates that additional flow can be pushed from σ to τ . Alternatively, the lack of an augmenting path indicates that a flow is maximum. A resource-level increment $\Delta(A)$ for an event set $A \subseteq R_H$ is related to a flow in $F(R_H)$ as follows. We define the producer weight in A as $c(A^+) = \sum_{e^+ \in A} c(e^+)$ and the consumer weight in A as $c(A^-) = \sum_{e^- \in A} c(e^-)$. We also define the *producer residual* in A for a flow f of $F(R_H)$ as $r(A^+) = c(A^+) - f(\{\sigma\}, A)$, i.e., the total residual capacity of the edges from σ to A , and the *consumer residual* in A as $r(A^-) = c(A^-) - f(A, \{\tau\})$.

Lemma 1: $\Delta(A) = r(A^+) - r(A^-) + f(A, A^c)$.

Proof: $\Delta(A) = c(A^+) - c(A^-) = (c(A^+) - f(\{\sigma\}, A)) - (c(A^-) - f(A, \{\tau\})) = r(A^+) - (c(A^-) - f(A, \{\tau\})) - f(A, A^c) = r(A^+) - r(A^-) + f(A, A^c)$. \square

We now focus on predecessor sets such as P_X .

Lemma 2: $f(P_X, P_X^c) \leq 0$. Moreover, $f(P_X, P_X^c) = 0$ if and only if $f(e_1, e_2) = 0$ for each $e_1 \in P_X^c$ and $e_2 \in P_X$.

Proof: From the definition of predecessor set there is no edge $e_2 \rightarrow e_1$ in $F(R_H)$ with $e_1 \in P_X^c$ and $e_2 \in P_X$. Therefore, $f(e_2, e_1) \leq 0$ and $f(P_X, P_X^c) \leq 0$. The second condition can be demonstrated by observing that the sum of any number of non-positive numbers is 0 if and only if each number is 0. \square

Corollary 1: $\Delta(P_X) \leq r(P_X^+) - r(P_X^-)$.

Proof: Immediate from Lemma 1 and Lemma 2.

5.2 Maximum flows and maximum resource-level increments

We can now find the maximum resource-level increment set $P_{\max} \subseteq R_H$. Since $P_{\max} = \emptyset$ may be true, the maximum resource-level increment is always non-negative. The computation of the set uses a maximum flow f_{\max} of $F(R_H)$. We indicate with $r_{\max}(A)$ the producer/consumer residual of A computed for f_{\max} . The following fundamental theorem holds.

Theorem 1: Given a partial plan R_H , consider the (possibly empty) set P_{\max} of events that are reachable from the source σ in the residual network of some f_{\max} of $F(R_H)$. P_{\max} is the predecessor set with maximum $\Delta(P_{\max}) \geq 0$.

Proof: Assume that $r_{\max}(e^+) = 0$ for each $e^+ \in F(R_H)$. In this case no event is reachable from σ in the residual network, thus $P_{\max} = \emptyset$ and $\Delta(P_{\max}) = 0$. From Corollary 1, for any predecessor set P_X it is $\Delta(P_X) \leq -r_{\max}(P_X^-) \leq 0 = \Delta(P_{\max})$ and therefore $\Delta(P_{\max})$ is maximum.

Assume now $r_{\max}(e^+) > 0$ for some $e^+ \in F(R_H)$, in which case P_{\max} is not empty. The following three properties hold.

1. P_{\max} is a predecessor set.

If not, there will be an event $e_2 \notin P_{\max}$ such that $|e_1 e_2| \leq 0$ for some event $e_1 \in P_{\max}$. From the definition of $\text{Aprec}(R)$, however, we know that there must be a path in $\text{Aprec}(R)$ from e_1 to e_2 . Since this path will be present in $F(R_H)$ with all links having infinite capacity, the path will also always be present in any residual network for any flow. Therefore there is a path in the residual network going from σ to e_1 (by definition of P_{\max}) and then to e_2 . Therefore, $e_2 \in P_{\max}$, which is a contradiction.

2. $r_{\max}(P_{\max}^-) = 0$.

If not, there will be an event $e^- \in P_{\max}$ such that $r_{\max}(e^-) > 0$. We can therefore build an augmenting path of $F(R_H)$ as follows: 1) a path from σ to e^- with positive residual capacity which exists by definition of P_{\max} ; and 2) an edge $e^- \rightarrow \tau$ with positive residual capacity $r_{\max}(e^-)$. The existence of the augmenting path means that f_{\max} is not a maximum flow, which is a contradiction.

3. $f_{\max}(P_{\max}, P_{\max}^c) = 0$

Since P_{\max} is a predecessor set, from the proof of Lemma 2 we know that $f_{\max}(P_{\max}, P_{\max}^c) \leq 0$. If $f_{\max}(P_{\max}, P_{\max}^c) < 0$, then there is a pair of events $e_1 \in P_{\max}$ and $e_2 \in P_{\max}^c$ such that $f_{\max}(e_1, e_2) < 0$. This means that the residual capacity from e_1 to e_2 is positive and therefore there is an edge $e_1 \rightarrow e_2$ in the residual network. But by definition of P_{\max} , this means that $e_2 \in P_{\max}$, which is a contradiction.

Applying the properties of P_{\max} to Lemma 1, $\Delta(P_{\max}) = r_{\max}(P_{\max}^+) - r_{\max}(P_{\max}^-) + f_{\max}(P_{\max}, P_{\max}^c) = r_{\max}(P_{\max}^+) > 0$.

To prove the maximality of P_{\max} , observe from Corollary 1 that a non-empty predecessor set P_X has $\Delta(P_X) > 0$ only if $r_{\max}(e^+) > 0$ for some $e^+ \in P_X$. It is easy to see that P_X is the set of events reachable in the residual graph from $P_X^+ \subseteq P_{\max}^+$ and that the properties at points 2 and 3 above also hold for P_X . Therefore, $\Delta(P_X) = r_{\max}(P_X^+) \leq r_{\max}(P_{\max}^+) = \Delta(P_{\max})$, which proves the maximality of P_{\max} . \square

The construction of P_{\max} discussed before does not guarantee its uniqueness since it depends on a specific maximum flow among potentially many for $F(R_H)$. The following theorem proves, however, that P_{\max} is indeed unique for all maximum flows of $F(R_H)$ and that it contains the minimum number of events among all predecessor sets with maximum positive resource-level increment.

Theorem 3: *The predecessor set P_{\max} with maximum resource-level increment $\Delta(P_{\max})$ and with minimum number of events is unique across all maximum flows of $F(R_H)$.*

Proof: Consider two maximum flows $f_{\max,j}$ and $f_{\max,k}$ among all maximum flows of $F(R_H)$ and assume that they produce two distinct maximum resource-level increment predecessor sets, $P_{\max,j}$ and $P_{\max,k}$. From the maximality of their increment, it must be $\Delta(P_{\max,k}) = \Delta(P_{\max,j}) = \Delta_{\max}$. We can rewrite one maximum predecessor set as $P_{\max,j} = P_{j \cap k} \cup P_{k-j}$ where $P_{j \cap k} = P_{\max,j} \cap P_{\max,k}$ and $P_{k-j} = P_{\max,k} - P_{\max,j}$. The hypothesis of distinction of $P_{\max,j}$ and $P_{\max,k}$ yields $P_{k-j} \neq \emptyset$.

First we observe that $\Delta(P_{j \cap k}) = \Delta_{\max}$. If not, $\Delta(P_{j-k}) > 0$ and $\Delta(P_{k-j}) > 0$. But this means that $P_{j \cup k} = P_{\max,j} \cup P_{\max,k}$ is a predecessor set such that $\Delta(P_{j \cup k}) > \Delta_{\max}$, which is a contradiction. Consider now $P_{\max,j}$ and let us call $r_j(e)$ the residual $r_{\max}(e)$ computed in flow $f_{\max,j}$. Since $\Delta(P_{j \cap k}) = \Delta_{\max}$, it must be $r_j(e^+) = 0$ if $e^+ \in P_{j-k}$. Also, $r_j(e^-) = 0$ for each $e^- \in P_{\max,j}$. From Lemma 1, $\Delta(P_{j-k}) = f_{\max,j}(P_{j-k}, P_{j-k}^c) = 0$. From Lemma 2 it follows that $f_{\max,j}(P_{j-k}, P_{j \cap k}) = 0$. Hence, there cannot be a link in the residual network from an event in $P_{j \cap k}$ to one in P_{j-k} . Therefore, $e \in P_{j-k}$ is not reachable from σ in the residual network and $P_{j-k} = \emptyset$. Since this is true for any pair $\langle j, k \rangle$, P_{\max} is unique.

The same argument applied to P_{\max} and $P^\emptyset \subseteq P_{\max}$ proves the minimality of P_{\max} , where P^\emptyset is a predecessor set such that $\Delta(P^\emptyset) = \Delta_{\max}$. \square

6 Building Resource Envelopes

So far we know that the resource level for a schedule s at time $t \in H$ is equal to $L_s(t) = \Delta(C_H) + \Delta(P_X)$ for some predecessor set P_X . However, it is not immediately obvious that the converse also applies. Given any predecessor set P_X , we want to be able to determine a time $t_X \in H$, the *separation time*, and a schedule s_X , the *separation schedule*, such that all and only the events in $C_H \cup P_X$ are scheduled at or before time t_X . The existence of a separation schedule and a separation time is not obvious because of the upper-bound constraints in the STN, i.e., the metric links between events that do not contribute to the construction of $\text{Aprec}(R)$. If some event occurs too early with respect to t_X , an upper-bound constraint may force some event to occur before time t_X even if it is not a successor in $\text{Aprec}(R)$. We now show that indeed we can find a separation time and schedule for *any* P_X and therefore also for P_{\max} . For the latter, we show that t_X is one of the times at which the resource level is maximum over H for any schedule. This yields the maximum resource envelope L_{\max} if $H = [t, t]$ and we scan t over the horizon T .

6.1 Latest events

First we find the events in P_X that will be scheduled at time t_X . We say that e is a *latest event* of P_X if it is not a strict predecessor of any other event in P_X , i.e., for any $e' \in P_X$, $|e' - e| \geq 0$. We will call $P_{X,\text{late}}$ the set of all latest events in P_X . Also, we define $P_{X,\text{early}} = P_X - P_{X,\text{late}}$.

The following property holds between events in $P_{X,\text{late}}$ and $P_{X,\text{early}}$.

Property 1: *Any event $e_1 \in P_{X,\text{early}}$ is a strict predecessor of some event $e_2 \in P_{X,\text{late}}$ i.e., $|e_2 - e_1| < 0$.*

Proof: Since $e_1 \in P_{X, \text{early}}$, there must be an event $e_{11} \in P_X$ such that $|e_{11}e_1| < 0$. If $e_{11} \in P_{X, \text{late}}$, the property is proven. Otherwise, we can find a finite chain of events $e_2 \rightarrow e_{1k} \rightarrow \dots \rightarrow e_{11} \rightarrow e_1$ with $e_2 \in P_{X, \text{late}}$ and $|e_2e_{1k}| < 0$, $|e_{1k}e_{1j-1}| < 0$ and $|e_{11}e_1| < 0$, yielding $|e_2e_1| < 0$ for the triangular inequality of the shortest paths. If we could not find an $e_2 \in P_{X, \text{late}}$ to start such a finite chain, the chain would have to become a cycle of events in $P_{X, \text{early}}$, which contradicts the temporal consistency of R . \square

6.2 Separation Time for Latest Events

We can construct a separation time t_X at which we will schedule all latest events.

Lemma 3: *There is a time interval $[t_{X, \text{min}}, t_{X, \text{max}}]$ that intersects all time bounds $[et(e), lt(e)]$ with $e \in P_{X, \text{late}}$ and such that $\text{start}(H) \leq t_{X, \text{max}}$.*

Proof: There must be a time value in common among all time bounds in $P_{X, \text{late}}$. If not, there would be two events $e_1, e_2 \in P_{X, \text{late}}$ such that $et(e_1) > lt(e_2)$ and, from the triangular inequality, $|e_1e_2| \leq -et(e_1) + lt(e_2) < 0$, which is inconsistent with the definition of $P_{X, \text{late}}$. Observe that there must be an event $e \in P_{X, \text{late}}$ such that $lt(e) = t_{X, \text{max}}$. If $\text{start}(H) > t_{X, \text{max}}$, then $lt(e) < \text{start}(H)$, which contradicts $e \in R_H$. \square

We define the separation time as $t_X = \max(\text{start}(H), t_{X, \text{min}})$, with $t_X = \text{start}(H)$ if $P_X = \emptyset$. We can then show that each event in P_X^c can be scheduled after t_X .

Lemma 4: *For any event $e \in P_X^c$, $lt(e) > t_X$.*

Proof: By definition of R_H it must be $lt(e) > \text{start}(H)$. So we only need to consider the case in which $t_X = t_{X, \text{min}} > \text{start}(H)$. In this case there is at least one event $e_1 \in P_{X, \text{late}}$ such that $et(e_1) = t_{X, \text{min}}$. For this event it is $|e_1e| \leq -et(e_1) + lt(e)$. Since $e \in P_X^c$, it must be that $|e_1e| > 0$, otherwise e would follow in $\text{Apred}(R)$ an event in P_X . Therefore, $lt(e) \geq et(e_1) + |e_1e| > et(e_1) = t_{X, \text{min}}$. \square

6.3 Separation schedule for predecessors

We now build the separation schedule s_X for P_X and t_X , i.e., a schedule such that $s_X(e) \leq t_X$ for $e \in C_H \cup P_X$ and $s_X(e) > t_X$ for $e \in P_X^c \cup O_X$. Note that the following discussion holds also if $P_X = \emptyset$ and $t_X = \text{start}(H)$.

The following algorithm builds the separation schedule.

1. Schedule all $e \in P_{X, \text{late}}$ at t_X , i.e., $s_X(e) = t_X$.
2. Propagate time through R obtaining new time bounds $[et'(e), lt'(e)]$ for each $e \in E(R)$.
3. Schedule all events $e \in E(R) - P_{X, \text{late}}$ at their new latest time, i.e., $s_X(e) = lt'(e)$.

For s_X to be a schedule, it must be consistent with respect to R . We see that step 1 is consistent since: 1) t_X belongs to the intersection of all latest event time bounds; and 2) since for any pair of latest events $|e_1e_2| \geq 0$, scheduling one at t_X does not prevent any other latest events to be also scheduled at time t_X . Step 3 above is also consistent because it is always possible to schedule all events at their latest times without temporal repropagation.

Now we need to show that the property defining a separation schedule is satisfied for s_X . Note that we already know that it is satisfied for events in $P_{X, \text{late}}$. By definition it is also satisfied for events in C_H and O_H . Therefore, we need to show that it is satisfied for $P_{X, \text{early}}$ and P_X^c .

a) $lt'(e) \leq t_X$ for all $e \in P_{X, \text{early}}$

According to Property 1 we can pick an event $e_1 \in P_{X, \text{late}}$ such that $|e_1e| < 0$. From the triangular inequality we have $lt'(e) \leq lt'(e_1) + |e_1e| < lt'(e_1) = t_X$.

b) $lt'(e) > t_x$ for all $e \in P_x^c$.

From Lemma 4 we know that before temporal repropagation it was $lt(e) > t_x$. After it, either $lt'(e) = lt(e)$, in which case the condition is satisfied, or $lt'(e)$ has changed due to a propagation that starts from some event $e_1 \in P_{x, \text{late}}$. So it must be $lt'(e) = t_x + |e_1 e|$. Since $e \in P_x^c$, it must be $|e_1 e| > 0$, otherwise e would follow in $\text{Apred}(R)$ an event in P_x . Hence, $lt'(e) > t_x$.

We can now compute the maximum resource level for any schedule within the interval H . In the following, we indicate with $P_{\max}(R_H)$ the P_{\max} computed over $F(R_H)$.

Theorem 4: *The maximum resource level for any schedule of R over an interval $H \subseteq T$ is given by $\Delta(C_H) + \Delta(P_{\max}(R_H))$.*

Proof: We know that at any time $t \in H$ the events in R_H that are scheduled before t are a predecessor set P_x . For the resource level at time t it is always $\Delta(C_H) + \Delta(P_x) \leq \Delta(C_H) + \Delta(P_{\max}(R_H))$, the latter being the resource level at the separation time t_x for the separation schedule s_x . \square

There are two interesting special cases of Theorem 4.

Corollary 2: *The maximum possible resource consumption for R over T is equal to $\Delta(P_{\max}(R_T))$.*

This means that estimating the maximum possible resource consumption for a flexible plan over the entire time horizon has the same complexity as a maximum flow problem.

Corollary 3: $L_{\max}(t) = \Delta(C_t) + \Delta(P_{\max}(R_t))$.

The last formula tells us how to compute the resource-level envelope at a specific time. We now need to find an efficient algorithm to compute the resource-level envelope over the entire horizon T .

7 Efficient Computation of Resource Envelopes

From Corollary 3, the naïve approach to compute a resource-level envelope would be to iterate over all possible $t \in T$. However, we only need to compute L_{\max} at times when either C_t or R_t changes. This can only happen at $et(e)$ or $lt(e)$ for any $e \in E(R)$. Therefore we need to compute new levels for L_{\max} only $2N$ times, where N is the number of start/end events in the original activity network. For each such computation, we need to: a) compute $P_{\max}(R_t)$ by running a maximum flow on a network with at most N nodes; and 2) collect and sum the events in C_t and $P_{\max}(R_t)$. The total complexity of the algorithm is therefore $O(N \cdot O(\text{maxflow}(N)) + N^2)$, where $O(\text{maxflow}(N))$ is the complexity of finding a maximum flow with an arbitrary maximum flow algorithm. For modern algorithms using the “preflow push” method [9], the worst case complexity can be $O(N^3)$. Extensive empirical studies show that the practical complexity of variations of the method can be as fast as $O(N^{1.5})$ [1]. This suggests that resource-level envelopes could operate in the inner loop of scheduling search algorithms, especially if they can be computed incrementally.

8 Conclusions

In this paper we describe an efficient algorithm to compute the tightest exact bound on the resource level induced by a flexible activity plan. This can potentially save exponential

amounts of work with respect to currently available looser bounds. Future work will pursue two directions. The first is developing more incremental algorithms for the computation of the envelope. Using a temporal scanning of the events in the temporal network, it should be possible to significantly reduce the size of the networks on which the maximum flow algorithm needs to be run. This could significantly speed up the envelope calculation. The second direction will test the practical effectiveness of resource envelopes in the inner loop of search algorithms for multi-capacity resource scheduling, such as those used in (Laborie, 2001). This includes inner-loop backtracking and termination tests and variable and value ordering heuristics that exploit more directly the properties of the resource envelopes.

Acknowledgements

Ari Jonsson and Jeremy Frank were instrumental in pushing me to focus on this problem. During a dinner discussion at possibly the worst tourist restaurant in Paris, Grigore Rosu convinced me that the key of the resource-level envelope problem lies with the maximum-flow problem. Paul Morris gave me several helpful comments and suggested a simplification of the proof of theorem 1. Finally, Amedeo Cesta, Mary Bernardine Dias, Gregory Dorais, Paul Tompkins, an anonymous reviewer of a previous, unsuccessful submission, and a reviewer of the current successful one, gave several comments that helped me improve the presentation. This work was performed with the support of the Intelligent Systems project of the Computing, Information and Communication Technologies research program of the National Aeronautics and Space Administration.

References

1. R.K. Ahuja, M. Kodialam, A.K. Mishra, J.B. Orlin. Computational Investigations of Maximum Flow Algorithms. *European Journal of Operational Research*, Vol 97(3), 1997.
2. K.R. Baker. *Introduction to Sequencing and Scheduling*. Wiley, New York, 1974.
3. J.C. Beck, A.J. Davenport, E.D. Davis, M.S. Fox. Beyond Contention: Extending Texture-Based Scheduling Heuristics. in *Proceedings of AAAI 1997*, Providence, RI, 1997.
4. A., Cesta, A. Oddi, S.F. Smith, A Constraint-Based Method for Resource Constrained Project Scheduling with Time Windows, CMU RI Technical Report, February 2000.
5. A. Cesta, C. Stella. A time and Resource Problem for Planning Architectures. *Proceedings of the 4th European Conference on Planning (ECP 97)*. Toulouse, France, 1997.
6. H. S.F. Cooper Jr., The Loneliness of the Long-Duration Astronaut, *Air & Space/Smithsonian*, June/July 1996, available at <http://www.airspacemag.com/ASM/Mag/Index/1996/JJ/lda.html>
7. T.H. Cormen, C.E. Leiserson, R.L. Rivest. *Introduction to Algorithms*. Cambridge, MA, 1990.
8. R. Dechter, I. Meiri, J. Pearl. Temporal Constraint Networks. *Artificial Intelligence*, 49:61-95, May 1991.
9. A.V. Goldberg, R.E. Tarjan. A New Approach to the Maximum-Flow Problem. *Journal of the ACM*, Vol. 35(4), 1988.
10. P. Laborie, Algorithms for Propagating Resource Constraints in AI Planning and Scheduling: Existing Approaches and New Results, *Proceedings of ECP 2001*, Toledo, Spain, 2001.
11. P. Morris, N. Muscettola, T. Vidal. Dynamic Control of Plans with Temporal Uncertainty, in *Proceedings of IJCAI 2001*, Seattle, WA, 2001.
12. N. Muscettola. On the Utility of Bottleneck Reasoning for Scheduling. in *Proceedings of AAAI 1994*, Seattle, WA, 1994.
13. W.P.M. Nuijten. Time and Resource Constrained Scheduling: a Constraint Satisfaction Approach. PhD Thesis, Eindhoven University of Technology, 1994.
14. N. Sadeh. Look-ahead techniques for micro-opportunistic job-shop scheduling. PhD Thesis, Carnegie Mellon University, CMU-CS-91-102, 1991.
15. M. Zweben, M.S. Fox. *Intelligent Scheduling*. Morgan Kaufmann, San Francisco, 1994.

