



**HAL**  
open science

## On the complexity of bounded time and precision reachability for piecewise affine systems

Hugo Bazille, Olivier Bournez, Walid Gomaa, Amaury Pouly

► **To cite this version:**

Hugo Bazille, Olivier Bournez, Walid Gomaa, Amaury Pouly. On the complexity of bounded time and precision reachability for piecewise affine systems. *Theoretical Computer Science*, 2018, 735, pp.132-146. 10.1016/j.tcs.2016.09.021 . hal-04303982

**HAL Id: hal-04303982**

**<https://cnrs.hal.science/hal-04303982>**

Submitted on 5 Feb 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# On The Complexity of Bounded Time and Precision Reachability for Piecewise Affine Systems\*

Hugo Bazille<sup>3</sup>, Olivier Bournez<sup>1</sup>, Walid Gomaa<sup>2,4</sup>, and Amaury Pouly<sup>1</sup>

<sup>1</sup> École Polytechnique, LIX, 91128 Palaiseau Cedex, France

<sup>2</sup> Egypt Japan University of Science and Technology, CSE, Alexandria, Egypt

<sup>3</sup> ENS Cachan/Bretagne et Université Rennes 1, France

<sup>4</sup> Faculty of Engineering, Alexandria University, Alexandria, Egypt

**Abstract.** Reachability for piecewise affine systems is known to be undecidable, starting from dimension 2. In this paper we investigate the exact complexity of several decidable variants of reachability and control questions for piecewise affine systems. We show in particular that the region-to-region bounded time versions leads to NP-complete or co-NP-complete problems, starting from dimension 2. We also prove that a bounded precision version leads to *PSPACE*-complete problems.

## 1 Introduction

A (discrete time) dynamical system  $\mathcal{H}$  is given by some space  $X$  and a function  $f : X \rightarrow X$ . A trajectory of the system starting from  $x_0$  is a sequence  $x_0, x_1, x_2, \dots$  etc., with  $x_{i+1} = f(x_i) = f^{[i+1]}(x_0)$  where  $f^{[i]}$  stands for  $i^{\text{th}}$  iterate of  $f$ . A crucial problem in such systems is the *reachability question*: given a system  $\mathcal{H}$  and  $R_0, R \subseteq X$ , determine if there is a trajectory starting from a point of  $R_0$  that falls in  $R$ . Reachability is known to be *undecidable* for very simple functions  $f$ . Indeed, it is well-known that various types of dynamical systems, such as hybrid systems, piecewise affine systems, or saturated linear systems, can simulate Turing machines, see e.g., [12,9,14,15].

This question is at the heart of *verification* of systems. Indeed, a safety property corresponds to the determination if there is a trajectory starting from some set  $R_0$  of possible initial states to the set  $R$  of bad states. The industrial and economical impact of having efficient computer tools, that are able to guarantee that a given system does satisfy its specification, have indeed generated very important literature. Particularly, many undecidability and complexity-theoretic results about the hardness of verification of safety properties have been obtained in the model checking community. However, as far as we know, the exact complexity of *natural restrictions* of the reachability question for systems as simple as piecewise affine maps are not known, despite their practical interest.

Indeed, existing results mainly focus on the frontier between decidability and undecidability. For example, it is known that reachability is undecidable

---

\* This work was partially supported by DGA Project CALCULS.

for piecewise constant derivative systems of dimension 3, whereas it is decidable for dimension 2 [1]. It is known that piecewise affine maps of dimension 2 can simulate Turing machines [13], whereas the question for dimension 1 is still open and can be related to other natural problems [2,3,5]. Variations of such problems over the integers have recently been investigated [6].

Some complexity facts follow immediately from these (un)computability results: for example, point-to-point bounded time reachability for piecewise affine maps is  $P$ -complete as it corresponds to configuration to configuration reachability for Turing machines.

However, there remain many natural variants of reachability questions for whose complexity have not yet been established.

For example, in the context of verification, proving the safety of a system can often be reduced to a reachability question of the form point-to-region or region-to-region reachability. These variants are more general questions than point-to-point reachability. Their complexities do not follow from existing results.

In this paper we choose to restrict to the case of piecewise affine maps and we consider the following natural variant of the problem.

**BOUNDED TIME:** we want to know if region  $R$  is reached in less than some prescribed time  $T$ , with  $f$  assumed to be continuous.

**FIXED PRECISION:** given some precision  $\varepsilon = 2^{-n}$ , we want to know if region  $R$  is reached by  $\tilde{f}_\varepsilon$  such that  $\tilde{f}_\varepsilon(x) = \lfloor \frac{f(x)}{\varepsilon} \rfloor \varepsilon$ . This corresponds to truncating  $f$  at precision  $2^{-n}$  on each coordinate, or equivalently assuming that computations happen at some precision not more than  $2^{-n}$ , for some given  $n$ . In other words, one wants to know if  $R$  is reached by the dynamics where precision operation is applied at each iteration.

*Remark 1.* We consider a version where everything is rounded downwards to a multiple of epsilon. Variants could also be considered, with for example closest or upper rounding. This would not change the complexity.

*Remark 2.* A variant could also be *chain reachability*: Deciding the existence of sequence  $x_n$  from initial region to target region such that at any intermediate step  $i$ ,  $\|x_{i+1} - f(x_i)\| \leq \epsilon$ . We do not know the complexity of variants based on this idea.

*Remark 3.* We consider piecewise affine maps over the domain  $[0, 1]^d$ , The case of integer domains has been studied in [6] and turns out to be quite different. We also assume  $f$  to be continuous. This makes the hardness result more natural.

In an orthogonal way, control of systems or constructions of controllers for systems often yield to dual questions. Instead of asking if some trajectory reaches region  $R$ , one wants to know if all trajectories reach  $R$ . The questions of stability, mortality, or nilpotence for piecewise affine maps and saturated linear systems have been established in [7]. Still in this context, the complexity of the problem when restricting to bounded time or fixed precision is not known.

This paper provides an exact characterization of the *algorithmic complexity* of those two types of reachability for discrete time dynamical systems. Let  $PAF_d$  denote the set of piecewise-affine *continuous* functions over  $[0, 1]^d$ . At the end we get the following picture.

*Remark 4.* Notice that we expect in several statements time to be given in unary, in order to get a completeness result. We do not know about the complexity of the variants where time would be given in binary.

**Problem:** REACH-REGION

**Inputs:** a continuous  $PAF_d$   $f$  and two regions  $R_0$  and  $R$  in  $\text{dom}(f)$

**Output:**  $\exists x_0 \in R_0, t \in \mathbb{N}, f^{[t]}(x_0) \in R?$

**Theorem 5 ([13]).** *Problem REACH-REGION is undecidable (and is recursively enumerable-complete).*

**Problem:** CONTROL-REGION

**Inputs:** a continuous  $PAF_d$   $f$  and two regions  $R_0$  and  $R$  in  $\text{dom}(f)$

**Output:**  $\forall x_0 \in R_0, \exists t \in \mathbb{N}, f^{[t]}(x_0) \in R?$

**Theorem 6 ([7]).** *Problem CONTROL-REGION is undecidable (and is co-recursively enumerable complete) for  $d \geq 2$ .*

**Problem:** REACH-REGION-TIME

**Inputs:** a time  $T \in \mathbb{N}$  in unary, a continuous  $PAF_d$   $f$  and two regions  $R_0$  and  $R$  in  $\text{dom}(f)$

**Output:**  $\exists x_0 \in R_0, \exists t \leq T, f^{[t]}(x_0) \in R?$

**Theorem 7 (Theorems 31 and 36).** *REACH-REGION-TIME is NP-complete for  $d \geq 2$ .*

**Problem:** CONTROL-REGION-TIME

**Inputs:** a time  $T \in \mathbb{N}$  in unary, a continuous  $PAF_d$   $f$  and two regions  $R_0$  and  $R$  in  $\text{dom}(f)$

**Output:**  $\forall x_0 \in R_0, \exists t \leq T, f^{[t]}(x_0) \in R?$

**Theorem 8 (Theorems 37 and 38).** *CONTROL-REGION-TIME is coNP-complete for  $d \geq 2$ .*

**Problem:** REACH-REGION-PRECISION

**Inputs:** a continuous  $PAF_d$   $f$ , two regions  $R_0$  and  $R$  in  $\text{dom}(f)$  and  $\varepsilon = 2^{-n}$ ,  $n$  given in unary

**Output:**  $\exists x_0 \in R_0, t \in \mathbb{N}, (\tilde{f}_\varepsilon)^{[t]}(x_0) \in R?$

**Problem:** CONTROL-REGION-PRECISION

**Inputs:** a continuous  $PAF_d$   $f$ , two regions  $R_0$  and  $R$  in  $\text{dom}(f)$  and  $\varepsilon = 2^{-n}$ ,  $n$  given in unary

**Output:**  $\forall x_0 \in R_0, \exists t \in \mathbb{N}, (\tilde{f}_\varepsilon)^{[t]}(x_0) \in R?$

**Theorem 9 (Theorems 40 and 40).** *REACH-REGION-PRECISION is PSPACE-complete for  $d \geq 2$ .*

**Theorem 10 (Theorems 40 and 40).** *CONTROL-REGION-PRECISION is PSPACE-complete for  $d \geq 2$ .*

All our problems are region-to-region reachability questions, and requires new proof techniques.

Indeed, classical tricks to simulate a Turing machine using a piecewise affine maps encode a Turing machine configuration by a point, and assume that all the points of the trajectories encode (possibly ultimately) valid Turing machines configurations.

This is not a problem in the context of point-to-point reachability, but this can not be extended to region-to-region reachability. Indeed, a (non-trivial) region consists mostly in invalid points: almost all points do not correspond to encoding of Turing machines for all the encodings considered in [13,7].

In order to establish hardness results, the trajectories of all (valid and invalid) points must be carefully controlled. This turns out not to be easily possible using the classical encodings.

Let us insist on the fact that we restrict our results to continuous dynamics. In this context, this is an additional source of difficulties: Dealing with points and trajectories not corresponding to valid configurations or evolutions.

A short version of this paper has been presented at the conference “Reachability Problems 2014” [4]. The current journal version contains full proofs for all statements, and is also providing new results: bounded precision variants (problems *REACH-REGION-PRECISION* and *CONTROL-REGION-PRECISION*) were not considered in short version [4].

## 2 Preliminaries

### 2.1 Notations

The set of non-negative integers is denoted  $\mathbb{N}$  and the set of the first  $n$  naturals is denoted  $\mathbb{N}_n = \{0, 1, \dots, n - 1\}$ . For any finite set  $\Sigma$ , let  $\Sigma^*$  denote the set of finite words over  $\Sigma$ . For any word  $w \in \Sigma^*$ , let  $|w|$  denote the length of  $w$ . Finally, let  $\lambda$  denote the empty word. If  $w$  is a word, let  $w_1$  denote its first character,  $w_2$  the second one and so on. For any  $i, j \in \mathbb{N}$ , let  $w_{i\dots j}$  denote the subword  $w_i w_{i+1} \dots w_j$ . For any  $\sigma \in \Sigma$ , and  $k \in \mathbb{N}$ , let  $\sigma^k$  denote the word of length  $k$  where all symbols are  $\sigma$ . For any function  $f$ , let  $f \upharpoonright E$  denote the restriction of  $f$  to  $E$  and let  $\text{dom}(f)$  denote the domain of definition of  $f$ .

### 2.2 Piecewise affine functions

Let  $I$  denote the unit interval  $[0, 1]$ . Let  $d \in \mathbb{N}$ . A convex closed polyhedron in the space  $I^d$  is the solution set of some linear system of inequalities:

$$Ax \leq \mathbf{b} \tag{1}$$

with coefficient matrix  $A$  and offset vector  $\mathbf{b}$ . Let  $PAF_d$  denote the set of piecewise-affine continuous functions over  $I^d$ : That is to say, any  $f: I^d \rightarrow I^d$  in  $PAF_d$ ,  $f$  satisfies:

- $f$  is continuous,
- there exists a sequence  $(P_i)_{1 \leq i \leq p}$  of convex closed polyhedra with nonempty interior such that  $f_i = f \upharpoonright P_i$  is affine,  $I^d = \bigcup_{i=1}^p P_i$  and  $\overset{\circ}{P}_i \cap \overset{\circ}{P}_j = \emptyset$  for  $i \neq j$ , where  $\overset{\circ}{P}$  denotes the interior of  $P$ .

In the following discussion we will always assume that any polyhedron  $P$  can be defined by a finite set of linear inequalities, where all the elements of  $A$  and  $\mathbf{b}$  in (1) are all rationals. A polyhedron over which  $f$  is affine will also be called a region.

### 2.3 Decision problems

In this paper, we will show hardness results by reduction from known hard problems. We give the statement of these latter problems in the following.

**Problem:** SUBSET-SUM

**Inputs:** a goal  $B \in \mathbb{N}$  and integers  $A_1, \dots, A_n \in \mathbb{N}$ .

**Output:**  $\exists I \subseteq \{1, \dots, n\}, \sum_{i \in I} A_i = B?$

**Theorem 11** ([8]). *SUBSET-SUM is NP-complete.*

**Problem:** NOSUBSET-SUM

**Inputs:** a witness  $B \in \mathbb{N}$  and integers  $A_1, \dots, A_n \in \mathbb{N}$ .

**Output:**  $\forall I \subseteq \{1, \dots, n\}, \sum_{i \in I} A_i \neq B?$

**Theorem 12.** *NOSUBSET-SUM is coNP-complete.*

*Proof.* Basically the same proof as Theorem 11 [8]. □

**Problem:** Linspace-WORD

**Inputs:** A Linear Bounded Automaton (i.e. a one-tape TM that does not use any space besides the input)  $\mathcal{M}$  and a word  $w \in \Sigma^*$ .

**Output:** does  $\mathcal{M}$  accept  $w$ ?

**Theorem 13** (see e.g. [10]). *Linspace-WORD is PSPACE-complete.*

## 3 Hardness of Bounded Time Reachability

In this section, we will show that REACH-REGION-TIME is an NP-hard problem by reducing SUBSET-SUM to it.

### 3.1 Solving SUBSET-SUM by iteration

We will now show how to solve the SUBSET-SUM problem by iterating a function. Consider an instance  $\mathcal{I} = (B, A_1, \dots, A_n)$  of SUBSET-SUM. We will need to introduce some notions before defining our piecewise affine function. Our first notion is that of configurations, representing partial summation of the number for a given choice of  $I$ .

*Remark 14.* Without loss of generality, we will only consider instances where  $A_i \leq B$ , for all  $i$ . Indeed, if  $A_i > B$ , it will never be an element of a solution to the instance and so we can simply remove this variable from the problem. This ensures that  $A_i < B + 1$  in everything that follows.

**Definition 15 (Configuration).** A configuration of  $\mathcal{I}$  is a tuple  $(i, \sigma, \varepsilon_i, \dots, \varepsilon_n)$  where  $i \in \{1, \dots, n+1\}$ ,  $\sigma \in \{0, \dots, B+1\}$ ,  $\varepsilon_j \in \{0, 1\}$  for all  $i \leq j$ . Let  $\mathcal{C}_{\mathcal{I}}$  be the set of all configurations of  $\mathcal{I}$ .

The intuitive understanding of a configuration, made formal in the next definition, is the following:  $(i, \sigma, \varepsilon_i, \dots, \varepsilon_n)$  represents a situation where after having summed a subset of  $\{A_1, \dots, A_{i-1}\}$ , we got a sum  $\sigma$  and  $\varepsilon_j$  is 1 if and only if we are to pick  $A_j$  in the future.

**Definition 16 (Transition function).** The transition function  $T_{\mathcal{I}} : \mathcal{C}_{\mathcal{I}} \rightarrow \mathcal{C}_{\mathcal{I}}$ , is defined as follows:

$$T_{\mathcal{I}}(i, \sigma, \varepsilon_i, \dots, \varepsilon_n) = \begin{cases} (i, \sigma) & \text{if } i = n+1 \\ (i+1, \min(B+1, \sigma + \varepsilon_i A_i), \varepsilon_{i+1}, \dots, \varepsilon_n) & \text{otherwise} \end{cases}$$

It should be clear, by definition of a subset sum that we have the following simulation result.

**Lemma 17.** For any configuration  $c = (i, \sigma, \varepsilon_i, \dots, \varepsilon_n)$  and  $k \in \{0, \dots, n+1-i\}$ ,

$$T_{\mathcal{I}}^{[k]}(c) = (i+k, \min(B+1, \sigma + \sum_{j=i}^{i+k-1} \varepsilon_j A_j), \varepsilon_{i+k}, \dots, \varepsilon_n)$$

*Proof.* By induction. □

A consequence of this simulation is that we can reformulate satisfiability in terms of reachability.

**Lemma 18.**  $\mathcal{I}$  is a satisfiable instance (i.e., admits a subset sum equal to the target value) if and only if there exists a configuration  $c = (1, 0, \varepsilon_1, \dots, \varepsilon_n) \in \mathcal{C}_{\mathcal{I}}$  such that  $T_{\mathcal{I}}^{[n]}(c) = (n+1, B)$ .

*Proof.* The “only if” direction is the simplest: assume there exists  $I \subseteq \{1, \dots, n\}$  such that  $\sum_{i \in I} A_i = B$ . Define  $\varepsilon_i = 1$  if  $i \in I$  and 0 otherwise. We get that  $\sum_{i=1}^n \varepsilon_i A_i = B$ . Apply Lemma 17 to get that:

$$\begin{aligned} T_{\mathcal{I}}^{[n]}(1, 0, \varepsilon_1, \dots, \varepsilon_n) &= (n+1, \min\left(B+1, 0 + \sum_{i=1}^n \varepsilon_i A_i\right)) \\ &= (n+1, \min(B+1, B)) = (n+1, B) \end{aligned}$$

The “if” direction is very similar: assume that there exists  $c = (1, 0, \varepsilon_1, \dots, \varepsilon_n)$  such that  $T_{\mathcal{I}}^{[n]}(c) = (n + 1, B)$ . Lemma 17 gives:

$$T_{\mathcal{I}}^{[n]}(1, 0, \varepsilon_1, \dots, \varepsilon_n) = (n + 1, \min \left( B + 1, 0 + \sum_{i=1}^n \varepsilon_i A_i \right))$$

We can easily conclude that  $\sum_{i=1}^n \varepsilon_i A_i = B$  and thus by defining  $I = \{i \mid \varepsilon_i = 1\}$  we get that  $\sum_{i \in I} A_i = B$ . Hence,  $\mathcal{I}$  is satisfiable.  $\square$

### 3.2 Solving SUBSET-SUM with a piecewise affine function

In this section, we explain how to simulate the function  $T_{\mathcal{I}}$  using a piecewise affine function and some encoding of the configurations for a given  $\mathcal{I} = (B, A_1, \dots, A_n)$ . Since the reduction is quite technical, we start with some intuitions. In order to simulate the function  $T_{\mathcal{I}}$ , we first need to encode configurations with real numbers. Let  $c = (i, \sigma, \varepsilon_1, \dots, \varepsilon_n)$  be a configuration, we encode it using two real numbers in  $[0, 1]$ : the first one encodes  $i$  and  $\sigma$  and the second one encodes  $\varepsilon_1, \dots, \varepsilon_n$ . A simple approach is to encode them as digits of dyadic numbers, as depicted below:

$$\langle c \rangle = \left( 0. \overbrace{0 \dots 0}^i \overbrace{0 \dots 0}^{\sigma} \varepsilon_1 \dots \varepsilon_n \right) = \left( \begin{array}{l} i2^{-p} + \sigma2^{-q} \\ \varepsilon_1 2^{-1} + \varepsilon_2 2^{-2} + \dots \end{array} \right).$$

In the above encoding, we allocate  $p$  bits to  $i$  and  $q - p$  bits to  $\sigma$ . We simply need to choose them large enough to accommodate the largest possible value. The rationale behind this encoding is that it is easy to implement the transition function  $f_{\mathcal{I}}$  with a linear function. Note that for technical reasons explained later, we need to encode the second coordinate in basis  $\beta = 5$  instead of 2. We now encode 0 as  $0^* = 1$  and 1 as  $1^* = 4$ . Graphically, the action of  $f$  is very simple:

$$\begin{aligned} \text{if } \varepsilon_i = 0^* \text{ then } f_{\mathcal{I}} \left( 0. \overbrace{0 \dots 0}^i \overbrace{0 \dots 0}^{\sigma} \varepsilon_i \varepsilon_{i+1} \dots \varepsilon_n \right) &= \left( 0. \overbrace{0 \dots 0}^{i+1} \overbrace{0 \dots 0}^{\sigma} \varepsilon_{i+1} \dots \varepsilon_n \right), \\ \text{if } \varepsilon_i = 1^* \text{ then } f_{\mathcal{I}} \left( 0. \overbrace{0 \dots 0}^i \overbrace{0 \dots 0}^{\sigma} \varepsilon_i \varepsilon_{i+1} \dots \varepsilon_n \right) &= \left( 0. \overbrace{0 \dots 0}^{i+1} \overbrace{0 \dots 0}^{\sigma + A_i} \varepsilon_{i+1} \dots \varepsilon_n \right). \end{aligned}$$

For technical reasons, we will in fact split the case  $\varepsilon_i = 1^*$  into two, depending on whether  $\sigma + A_i$  becomes greater than  $B + 1$  or not. This is similar to what we did in Definition 16 with  $\min(B + 1, \sigma + A_i)$ :

$$\text{if } \sigma + A_i > B \text{ then } f_{\mathcal{I}} \left( 0. \overbrace{0 \dots 0}^i \overbrace{0 \dots 0}^{\sigma} \varepsilon_i \varepsilon_{i+1} \dots \varepsilon_n \right) = \left( 0. \overbrace{0 \dots 0}^{i+1} \overbrace{0 \dots 0}^{B+1} \varepsilon_{i+1} \dots \varepsilon_n \right).$$



We can formulate the “if  $\varepsilon_i = \alpha$ ” by using regions. Indeed, the set of encodings such that  $\varepsilon_i = \alpha$  is

$$\left\{ \begin{pmatrix} 0. \begin{array}{|c|c|} \hline i & \sigma \\ \hline \end{array} \\ \hline 0. \begin{array}{|c|c|c|c|c|} \hline 0 & \cdots & \alpha & \varepsilon_{i+1} & \cdots & \varepsilon_n \\ \hline \end{array} \end{pmatrix} : \sigma \in \mathbb{N}, \varepsilon_j \in \{0, 1\} \right\} \subset \begin{array}{l} [i2^{-p}, i2^{-p} + 2^{-p-1}] \\ \times [\alpha\beta^{-i}, (\alpha + 1)\beta^{-i}] \end{array}.$$

It is crucial to note that the statement of the problems only allows for polyhedral regions. This is why in the above equation, we had to overapproximate the region by intervals on each coordinate. This overapproximation is the root of all difficulties. Indeed, the region now contains many points that do not correspond to encodings anymore. We now go to the details of the construction.

**Definition 19 (Encoding).** Define  $p = \lceil \log_2(n + 2) \rceil$ ,  $\omega = \lceil \log_2(B + 2) \rceil$ ,  $q = p + \omega + 1$  and  $\beta = 5$ . Also define  $0^* = 1$  and  $1^* = 4$ . For any configuration  $c = (i, \sigma, \varepsilon_i, \dots, \varepsilon_n)$ , define the encoding of  $c$  as follows:

$$\langle c \rangle = \left( i2^{-p} + \sigma 2^{-q}, 0^* \beta^{-n-1} + \sum_{j=i}^n \varepsilon_j^* \beta^{-j} \right)$$

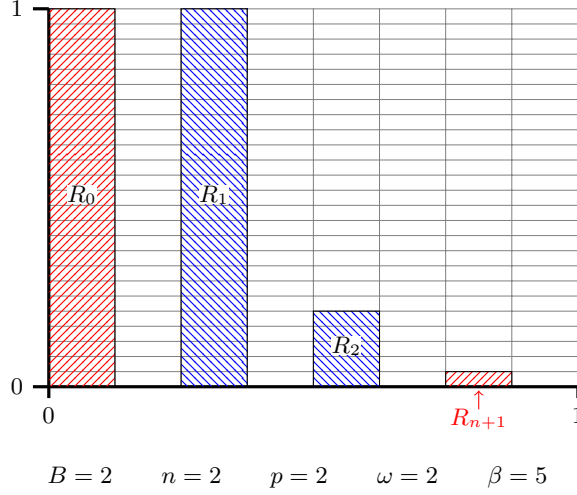
Also define the following regions for any  $i \in \{1, \dots, n+1\}$  and  $\alpha \in \{0, \dots, \beta-1\}$ :

$$\begin{aligned} R_0 &= [0, 2^{-p-1}] \times [0, 1] & R_i &= [i2^{-p}, i2^{-p} + 2^{-p-1}] \times [0, \beta^{-i+1}] \quad (i \geq 1) \\ R_{i,\alpha} &= [i2^{-p}, i2^{-p} + 2^{-p-1}] \times [\alpha\beta^{-i}, (\alpha + 1)\beta^{-i}] & R_i &= \cup_{\alpha \in \mathbb{N}_\beta} R_{i,\alpha} \\ R_{i,1^*}^{lin} &= [i2^{-p}, i2^{-p} + (B + 1 - A_i)2^{-q}] \times [1^* \beta^{-i}, 5\beta^{-i}] & R_{i,1^*}^{sat} &= R_{i,1^*} \setminus R_{i,1^*}^{lin} \end{aligned}$$

As noted before, we use basis  $\beta = 5$  on the second component to get some “space” between consecutive encodings. The choice of the value 1 and 4 for the encoding of 0 and 1, although not crucial, has been made to simplify the proof as much as possible.

The region  $R_0$  is for initialization purposes and is defined differently from the other  $R_i$ . The regions  $R_i$  correspond to the different values of  $i$  in the configuration (the current number). Each  $R_i$  is further divided into the  $R_{i,\alpha}$  corresponding to all the possible values of the next  $\varepsilon$  variable (recall that it is encoded in basis  $\beta$ ). In the special case of  $\varepsilon = 1$ , we cut the region  $R_{i,1^*}$  into a linear part and a saturated part. This is needed to emulate the  $\min(\sigma + A_i, B + 1)$  in Definition 16: the linear part corresponds to  $\sigma + A_i$  and the saturated part to  $B + 1$ . Figure 1 and Figure 2 give a graphical representation of the regions.

**Lemma 20.** For any configuration  $c = (i, \sigma, \varepsilon_i, \dots, \varepsilon_n)$ , if  $i = n + 1$  then  $\langle c \rangle \in R_{n+1,0^*}$ , otherwise  $\langle c \rangle \in R_{i,\varepsilon_i^*}$ . Furthermore if  $\varepsilon_i = 1$  and  $\sigma + A_i \leq B + 1$ , then  $\langle c \rangle \in R_{i,1^*}^{lin}$ , otherwise  $\langle c \rangle \in R_{i,1^*}^{sat}$ .



**Fig. 1.** Graphical representation of the regions

*Proof.* Recall that  $\omega = \lceil \log_2(B + 2) \rceil$  so  $B + 1 < 2^\omega$ , and  $q = p + \omega + 1$ . Since  $\sigma \leq B + 1$  by definition,  $(n + 1)2^{-p} \leq \langle c \rangle_1 \leq (n + 1)2^{-p} + (B + 1)2^{-p-1-\omega} \leq (n + 1)2^{-p} + 2^{-p-1}$ . This shows the result for the first component. In the case where  $\sigma + A_i \leq B + 1$  then  $\sigma 2^{-q} \leq (B + 1 - A_i)2^{-q}$  yielding the result for the second part of the result for the first component.

If  $i = n + 1$ , then  $\langle c \rangle_2 = 0^* \beta^{-p-1}$  and it trivially belongs to  $[0^* \beta^{-n-1}, (0^* + 1) \beta^{-n-1}]$ . Otherwise,

$$\begin{aligned}
 \varepsilon_i^* \beta^{-i} \leq \langle c \rangle_2 &\leq \varepsilon_i^* \beta^{-i} + \sum_{j=i+1}^{n+1} 1^* \beta^{-j} \leq \varepsilon_i^* \beta^{-i} + 1^* \beta^{-i-1} \frac{1 - \beta^{n-i}}{1 - \beta^{-1}} \\
 &\leq \varepsilon_i^* \beta^{-i} + 4\beta^{-i-1} \frac{\beta}{\beta - 1} \leq \varepsilon_i^* \beta^{-i} + \beta^{-i} \leq (\varepsilon_i^* + 1) \beta^{-i}.
 \end{aligned}$$

This shows the result when  $i < n + 1$ , for the second component of the result.  $\square$

We can now define a piecewise affine function that will mimic the behavior of  $T^{\mathcal{L}}$ . The region  $R_0$  is here to ensure that we start from a “clean” value on the first coordinate.

**Definition 21 (Piecewise affine simulation).**

$$f_{\mathcal{I}}(a, b) = \begin{cases} (2^{-p}, b) & \text{if } (a, b) \in R_0 \\ (a, b) & \text{if } (a, b) \in R_{n+1} \\ (a + 2^{-p}, b - 0^* \beta^{-i}) & \text{if } (a, b) \in R_{i,0^*} \\ (a + 2^{-p} + A_i 2^{-q}, b - 1^* \beta^{-i}) & \text{if } (a, b) \in R_{i,1^*}^{lin} \\ ((i + 1)2^{-p} + (B + 1)2^{-q}, b - 1^* \beta^{-i}) & \text{if } (a, b) \in R_{i,1^*}^{sat} \end{cases}$$

**Lemma 22 (Simulation is correct).** *For any configuration  $c \in \mathcal{C}_{\mathcal{I}}$ ,  $\langle T_{\mathcal{I}}(c) \rangle = f_{\mathcal{I}}(\langle c \rangle)$ .*

*Proof.* Let  $c = (i, \sigma, \varepsilon_i, \dots, \varepsilon_n)$ . There are several cases to consider: if  $i = n + 1$  then  $T_{\mathcal{I}}(c) = c$ , also by Lemma 20,  $\langle c \rangle \in R_{n+1,0^*}$ . Thus by definition of  $f$ ,  $f_{\mathcal{I}}(\langle c \rangle) = \langle c \rangle = \langle T(c) \rangle$  and this shows the result. If  $i < n + 1$ , we have three more cases to consider: the case where we don't take the value ( $\varepsilon_i = 0$ ) and the two cases where we take it ( $\varepsilon_i = 1$ ) with and without saturation.

- If  $\varepsilon_i = 0$  then  $T_{\mathcal{I}}(c) = (i + 1, \sigma, \varepsilon_{i+1}, \dots, \varepsilon_n)$ . On the other hand,  $\langle c \rangle = (a, b) = (i2^{-p} + \sigma 2^{-q}, 0^* \beta^{-i} + \sum_{j=i+1}^n \varepsilon_j \beta^{-j} + 0^* \beta^{-n-1})$ . By Lemma 20,  $\langle c \rangle \in R_{i,0^*}$  so by definition of  $f$ :

$$\begin{aligned} f_{\mathcal{I}}(\langle c \rangle) &= (a + 2^{-p}, b - 0^* \beta^{-i}) \\ &= ((i + 1)2^{-p} + \sigma 2^{-q}, \sum_{j=i+1}^n \varepsilon_j \beta^{-j} + 0^* \beta^{-n-1}) \\ &= \langle (i + 1, \sigma, \varepsilon_{i+1}, \dots, \varepsilon_n) \rangle = \langle T_{\mathcal{I}}(c) \rangle \end{aligned}$$

- If  $\varepsilon_i = 1$  and  $\sigma + A_i \leq B + 1$  then  $T_{\mathcal{I}}(c) = (i + 1, \sigma + A_i, \varepsilon_{i+1}, \dots, \varepsilon_n)$ . On the other hand,  $\langle c \rangle = (a, b) = (i2^{-p} + \sigma 2^{-q}, 1^* \beta^{-i} + \sum_{j=i+1}^n \varepsilon_j \beta^{-j} + 0^* \beta^{-n-1})$ . By Lemma 20,  $\langle c \rangle \in R_{i,1^*}^{lin}$  so by definition of  $f$ :

$$\begin{aligned} f_{\mathcal{I}}(\langle c \rangle) &= (a + 2^{-p} + A_i 2^{-q}, b - 1^* \beta^{-i}) \\ &= ((i + 1)2^{-p} + (\sigma + A_i)2^{-q}, \sum_{j=i+1}^n \varepsilon_j \beta^{-j} + 0^* \beta^{-n-1}) \\ &= \langle (i + 1, \sigma + A_i, \varepsilon_{i+1}, \dots, \varepsilon_n) \rangle = \langle T_{\mathcal{I}}(c) \rangle \end{aligned}$$

- If  $\varepsilon_i = 1$  and  $\sigma + A_i > B + 1$  then  $T_{\mathcal{I}}(c) = (i + 1, B + 1, \varepsilon_{i+1}, \dots, \varepsilon_n)$ . By Lemma 20,  $\langle c \rangle \in R_{i,1^*}^{sat}$  so by definition of  $f$ :

$$\begin{aligned} f_{\mathcal{I}}(\langle c \rangle) &= ((i + 1)2^{-p} + (B + 1)2^{-q}, b - 1^* \beta^{-i}) \\ &= \langle (i + 1, B + 1, \varepsilon_{i+1}, \dots, \varepsilon_n) \rangle = \langle T_{\mathcal{I}}(c) \rangle \end{aligned}$$

□

### 3.3 Making the simulation stable

In the previous section, that we have defined  $f_{\mathcal{I}}$  over a subset of the entire space and it is clear that this subspace is not stable in any way<sup>5</sup>. In order to match the definition of a piecewise affine function, we need to define  $f$  over the entire space or a stable subspace (containing the initial region). We follow this second approach and extend the definition of  $f$  on some more regions. More precisely, we need to define  $f$  over  $R_i = R_{i,0} \cup R_{i,1} \cup R_{i,2} \cup R_{i,3} \cup R_{i,4,1}$  and at the moment we have only defined  $f$  over  $R_{i,1} = R_{i,0^*}$  and  $R_{i,4} = R_{i,1^*}$ . Also note that

<sup>5</sup> For example  $R_{1,1} \subseteq f(R_0)$  but  $f$  is not defined over  $R_{1,1}$ .

$R_{i,4} = R_{i,4}^{lin} \cup R_{i,4}^{sat}$  and we define  $f$  separately on those two subregions. In order to correctly and continuously extend  $f$ , we will need to further split the region  $R_{i,3}$  into linear and saturated parts  $R_{i,3}^{lin}$  and  $R_{i,3}^{sat}$ : see Figure 2.

Before jumping into the technical details of the extension, we start with the intuition. First, it is crucial to understand that our main constraint is continuity: since we already defined  $f$  over  $R_{i,1}$  and  $R_{i,4}$ , our extension need to agree with  $f$  on the borders of those regions. Furthermore,  $f$  still needs to be affine, leaving us with little flexibility. Second, the behaviour of  $f$  on those regions needs to be carefully chosen. Indeed, as we mentioned before, we had to overapproximate regions in several places and our simulation now includes extra points. We do not want these extra points to have completely unpredictable trajectories, otherwise they might reach the final region by chance and break the reduction. Therefore, our strategy is to define  $f$  in such a way that its behaviour on those “wrong points” still has a valid interpretation in the original SUBSET-SUM problem. We detail this idea for the various region right after.

Let  $(a, b) \in R_{i,0} \cup R_{i,2} \cup R_{i,3}$ : intuitively, this point corresponds to a configuration  $c = (i, \sigma, \varepsilon_i, \dots, \varepsilon_n)$  where  $\varepsilon_i \notin \{0^*, 1^*\}$ . We know by construction that this point does not correspond to a proper trajectory so it is tempting to simply discard it. A very simple way of discarding point is to send them to a point  $x$  that is (i) stable by  $f$  ( $f(x) = x$ ) and (ii) not in the accepting region. That way we trap the trajectory of invalid points into a useless region of the space. Let us illustrate this on  $R_{i,0}$ : let  $(a, b) \in R_{i,0}$  and take  $f_{\mathcal{I}}(a, b) = (a, b)$ . It is now trivial that the point is stuck in  $R_{i,0}$ . Unfortunately,  $f_{\mathcal{I}}$  is not continuous anymore. Indeed, for  $(a, b) = (a, \beta^{-i}) \in R_{i,0} \cap R_{i,1}$  we have a discontinuity on the first coordinate. Indeed,  $f_1(a, b) = a$  on one side but  $f_1(a, b) = a + 2^{-p}$  on the other. A simple fix is to take  $f(a, b) = (a + 2^{-p}, 0)$ , this corresponds to:

$$\text{if } \varepsilon_i = 0 \text{ then } f_{\mathcal{I}} \left( \begin{array}{c} \begin{array}{|c|c|} \hline i & \sigma \\ \hline \end{array} \\ \hline \begin{array}{|c|c|c|c|c|} \hline 0 & \cdots & \varepsilon_i & \varepsilon_{i+1} & \cdots & \varepsilon_n \\ \hline \end{array} \end{array} \right) = \left( \begin{array}{c} \begin{array}{|c|c|} \hline i+1 & \sigma \\ \hline \end{array} \\ \hline \begin{array}{|c|c|c|c|c|} \hline 0 & \cdots & 0 & 0 & \cdots & 0 \\ \hline \end{array} \end{array} \right).$$

One can check that  $f$  is also continuous on the second coordinate. The reason why this choice is clever is because  $f(R_{i,0}) \subseteq R_{i+1,0}$ . Although the invalid points are not stable, they are now stuck in the bottom regions  $R_{j,0}$  for the rest of the simulation.

Unfortunately, this trick now does not work on  $R_{i,2}$  because of the continuity requirement on the second coordinate. Indeed, on  $R_{i,1} \cap R_{i,2}$  we have that  $f(a, b) = (a + 2^{-p}, \beta^{-i})$ . To understand what it means, imagine a configuration where  $\varepsilon_i = 2$  and all the remaining  $\varepsilon_j$  are 0, then its image by  $f$  corresponds to a configuration where all the remaining  $\varepsilon_j$  are 4 = 1\* plus an error. Indeed,  $\beta^{-i} = \sum_{n=i+1}^{\infty} 4\beta^{-n}$ . In other words, we have:

$$\text{if } \varepsilon_i = 2 \text{ then } f_{\mathcal{I}} \left( \begin{array}{c} \begin{array}{|c|c|} \hline i & \sigma \\ \hline \end{array} \\ \hline \begin{array}{|c|c|c|c|c|} \hline 0 & \cdots & \varepsilon_i & 0 & \cdots & 0 \\ \hline \end{array} \end{array} \right) = \left( \begin{array}{c} \begin{array}{|c|c|} \hline i+1 & \sigma \\ \hline \end{array} \\ \hline \begin{array}{|c|c|c|c|c|} \hline 0 & \cdots & 0 & 4 & \cdots & 4 & 4 & \cdots \\ \hline \end{array} \end{array} \right).$$

Furthermore, thinking about the future a bit, on  $R_{i,3} \cap R_{i,4}$  we have that  $f(a, b) = (a + 2^{-p}, 0)$ . In other words, somewhere on  $R_{i,2} \cup R_{i,3}$ , the second coordinate

has to go from  $\beta^{-i}$  to 0 in a continuous fashion. This is where the clever tricks comes in: we can continuously change the second coordinate from 1 to 0 such that the action of  $f$  looks like all  $\varepsilon_j$  were “flipped”: 0 is exchange with 4 and 1 with 2. To visualise how this possible, simply think about the configuration as having an infinite number of  $\varepsilon_j$  (although we use a finite amount of them) that gets turned into an infinite number of  $\mu_j$  where  $\mu_j = 4 - \varepsilon_j$ :

$$f_{\mathcal{I}} \left( \begin{array}{c} \left( \begin{array}{c} 0. \quad \begin{array}{|c|c|} \hline i+1 & \sigma \\ \hline \end{array} \\ \hline \end{array} \right) \\ \left( \begin{array}{c} 0 \quad \dots \quad 2 \quad \varepsilon_{i+1} \dots \varepsilon_n \quad \varepsilon_{n+1} \dots \end{array} \right) \end{array} \right) = \left( \begin{array}{c} \left( \begin{array}{c} 0. \quad \begin{array}{|c|c|} \hline i+1 & \sigma \\ \hline \end{array} \\ \hline \end{array} \right) \\ \left( \begin{array}{c} 0 \quad \dots \quad 0 \quad \mu_{i+1} \dots \mu_n \mu_{n+1} \dots \end{array} \right) \end{array} \right).$$

This might seem convoluted at first until one realises why this is helpful. Imagine an extended SUBSET-SUM simulation where we now have three actions instead of two:

- $\varepsilon_i = 0$ : go to next number,
- $\varepsilon_i = 1$ : add  $A_i$  and go to next number,
- $\varepsilon_i = 2$ : flip all remaining  $\varepsilon_j$  and go to next number.

Our simulation corresponds to this extended problem, which is still NP-complete. The remaining region  $R_{i,3}$  follows the same principle as  $R_{i,0}$ , with a slight complication because of the saturated sum. Figure 3 gives the interpretation of the definition of  $f$  on each subregion.

**Definition 23 (Extended region splitting).** For  $i \in \{1, \dots, n\}$  and  $\alpha \in \{0, \dots, \beta - 1\}$ , define:

$$R_{i,3}^{lin} = R_{i,3} \cap \left\{ (a, b) \mid b\beta^i - 3 \leq \frac{2^{-p-1} + i2^{-p} - a}{2^{-p-1} - (B+1 - A_i)2^{-q}} \right\} \quad R_{i,3}^{sat} = R_{i,3} \setminus R_{i,3}^{lin}$$

It should be clear by definition that  $R_{i,3}^{sat} = R_{i,3}^{lin} \cup R_{i,3}^{sat}$  and that the two subregions are disjoint except on the border.

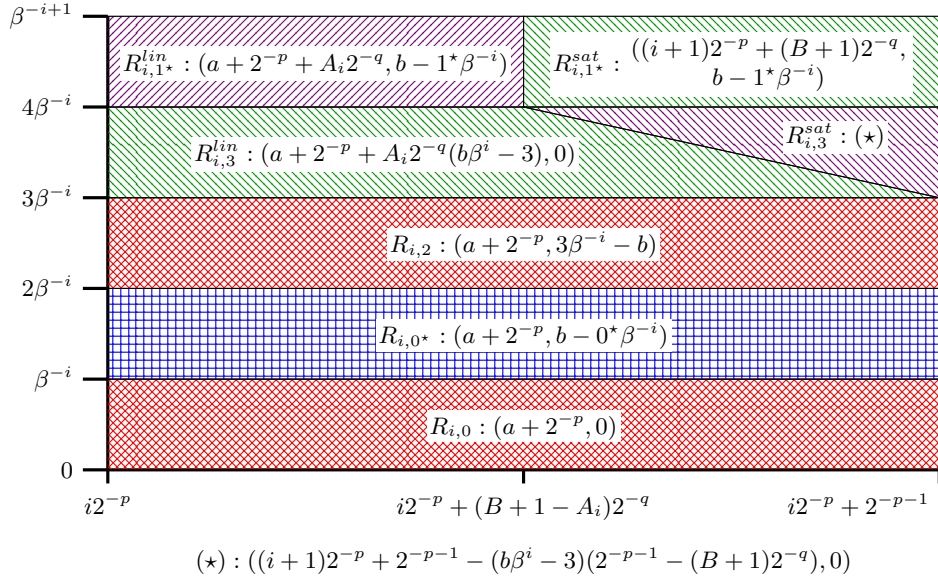
**Definition 24 (Extended piecewise affine simulation).**

$$f_{\mathcal{I}}(a, b) = \begin{cases} (a + 2^{-p}, 0) & \text{if } (a, b) \in R_{i,0} \\ (a + 2^{-p}, 3\beta^{-i} - b) & \text{if } (a, b) \in R_{i,2} \\ (a + 2^{-p} + A_i 2^{-q} (b\beta^i - 3), 0) & \text{if } (a, b) \in R_{i,3}^{lin} \\ \left( \left( i + \frac{3}{2} \right) 2^{-p} - (b\beta^i - 3)(2^{-p-1} - (B+1)2^{-q}), 0 \right) & \text{if } (a, b) \in R_{i,3}^{sat} \end{cases}$$

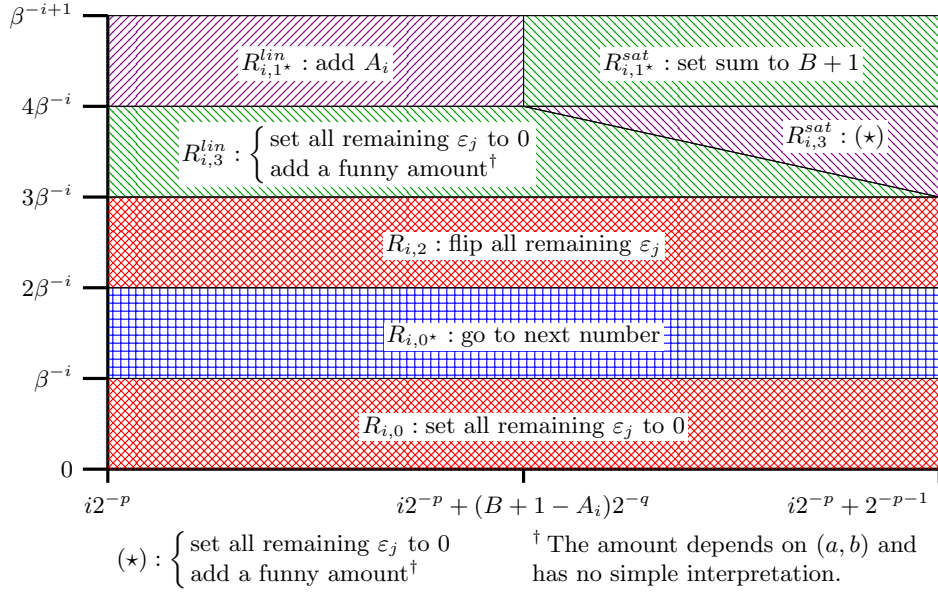
This extension was carefully chosen for its properties. In particular, we will see that  $f$  is still continuous, Also, the domain of definition of  $f$  is  $f$ -stable (i.e.  $f(\text{dom } f) \subseteq \text{dom } f$ ). And finally, we will see that  $f$  is somehow “reversible”.

**Lemma 25 (Simulation is continuous).** For any  $i \in \{1, \dots, n\}$ ,  $f_{\mathcal{I}}(R_i)$  is well-defined and continuous over  $R_i$ .

*Proof.* As outlined on Figure 2, we need to check that the definitions of  $f$  match at the borders of each subregions of  $R_i$ . More precisely, we need to check that Definition 21 and Definition 24 agree on all borders.



**Fig. 2.** Zoom on one  $R_i$  with the subregions and formulas.



**Fig. 3.** Interpretation of the behaviour of  $f$  on each subregion. See also Figure 2 and Section 3.2. All regions include an implicit “go to next number”.

- $(a, b) \in R_{i,0} \cap R_{i,0^*}$ : the first component is computed using the same formula so is clearly continuous, the second component is always 0 on both side of the border because  $b - 0^* \beta^{-i} = 0$  for  $b = \beta^{-i}$
- $(a, b) \in R_{i,0^*} \cap R_{i,2}$ : the first component is computed using the same formula so is clearly continuous, the second component is always  $\beta^{-i}$  on both side of the border because  $b - 0^* \beta^{-i} = 3\beta^{-i} - b = \beta^{-i}$  for  $b = 2\beta^{-i}$
- $(a, b) \in R_{i,2} \cap R_{i,3}^{lin}$ : the first component is  $a + 2^{-p}$  and the second component is 0 on both side of the border because  $3\beta^{-i} - b = b\beta^i - 3 = 0$  for  $b = 3\beta^{-i}$
- $(a, b) \in R_{i,3}^{lin} \cap R_{i,1^*}^{lin}$ : the first component is  $a + 2^{-p} + A_i 2^{-q}$  and the second component is 0 on both side of the border because  $b\beta^i - 3 = 1$  and  $b - 1^* \beta^{-i} = 0$  for  $b = 4\beta^{-i}$
- $(a, b) \in R_{i,3}^{lin} \cap R_{i,3}^{sat}$ : the second component is always 0 on both regions so is clearly continuous. From Definition 23 one can see that  $b\beta^i - 3 = \frac{Y}{X}$  holds on the border where  $Y = 2^{-p-1} + i2^{-p} - a$  and  $X = 2^{-p-1} - (B + 1 - A_i)2^{-q}$ . Consequently, if we compute the difference between the two expression at the borders, we get:

$$\begin{aligned}
& ((i+1)2^{-p} + 2^{-p-1} - (b\beta^i - 3)(2^{-p-1} - (B+1)2^{-q})) \\
& - (a + 2^{-p} + A_i 2^{-q}(b\beta^i - 3)) \\
& = i2^{-p} + 2^{-p-1} - a - \frac{Y}{X}(2^{-p-1} - (B+1)2^{-q} + A_i 2^{-q}) \\
& = Y - \frac{Y}{X}X = 0
\end{aligned}$$

this proves that they are equal.

- $(a, b) \in R_{i,3}^{sat} \cap R_{i,1^*}^{sat}$ : the first component is  $(i+1)2^{-p} + (B+1)2^{-q}$  and the second component is 0 on both side of the border because  $b\beta^i - 3 = 1$  and  $b - 1^* \beta^{-i} = 0$  for  $\beta = 4\beta^{-i}$
- $(a, b) \in R_{i,1^*}^{lin} \cap R_{i,1^*}^{sat}$ : the first component is  $(i+1)2^{-p} + (B+1)2^{-q}$  on both side of the border because  $a = i2^{-p} + (B+1 - A_i)2^{-q}$ , and the second component is computed using the same formula so is clearly continuous

□

**Lemma 26 (Simulation is stable).** *For any  $i \in \{1, \dots, n\}$ ,  $f_{\mathcal{I}}(R_i) \subseteq R_{i+1}$ . Furthermore,  $f(R_0) \subseteq R_1$  and  $f(R_{n+1}) \subseteq R_{n+1}$ .*

*Proof.* We need to examine all possible cases for  $(a, b) \in R_i$ . Since  $R_i = \bigcup_{\alpha=0}^4 R_{i,\alpha}$  and that  $R_{i,\alpha} = R_{i,\alpha}^{lin} \cup R_{i,\alpha}^{sat}$  we indeed cover all cases.

- If  $(a, b) \in R_0$ : then  $f_{\mathcal{I}}(a, b) = (a + 2^{-p}, b)$  so  $f_{\mathcal{I}}(R_0) = f_{\mathcal{I}}([0, 2^{-p-1}] \times [0, 1]) = [2^{-p}, 2^{-p} + 2^{-p-1}] \times [0, 1] = R_1$ .
- If  $(a, b) \in R_{n+1}$ : then  $f_{\mathcal{I}}(a, b) = (a, b)$  so  $f_{\mathcal{I}}(R_{n+1}) = R_{n+1}$ .
- If  $(a, b) \in R_{i,0}$ : then  $f_{\mathcal{I}}(a, b) = (a + 2^{-p}, 0)$  so  $f_{\mathcal{I}}(R_{i,0}) = f_{\mathcal{I}}([i2^{-p}, i2^{-p} + 2^{-p-1}] \times [0, \beta^{-i}]) = [(i+1)2^{-p}, (i+1)2^{-p} + 2^{-p-1}] \times \{0\} \subseteq R_{i+1}$ .
- If  $(a, b) \in R_{i,1} = R_{i,0^*}$ : then  $f_{\mathcal{I}}(a, b) = (a + 2^{-p}, b - 0^* \beta^{-i})$  so  $f_{\mathcal{I}}(R_{i,1}) = f_{\mathcal{I}}([i2^{-p}, i2^{-p} + 2^{-p-1}] \times [\beta^{-i}, 2\beta^{-i}]) = [(i+1)2^{-p}, (i+1)2^{-p} + 2^{-p-1}] \times [0, \beta^{-i}] = R_{i+1}$ .

- If  $(a, b) \in R_{i,2}$ : then  $f_{\mathcal{I}}(a, b) = (a+2^{-p}, 3\beta^{-i}-b)$  so  $f_{\mathcal{I}}(R_{i,2}) = f_{\mathcal{I}}([i2^{-p}, i2^{-p}+2^{-p-1}] \times [2\beta^{-i}, 3\beta^{-i}]) = [(i+1)2^{-p}, (i+1)2^{-p}+2^{-p-1}] \times [0, \beta^{-i}] = R_{i+1}$ .
- If  $(a, b) \in R_{i,3}^{lin}$ : the image of the second component is always 0 so it's easy for this one, also from Definition 23,  $b\beta^i - 3 \leq \frac{2^{-p-1}+i2^{-p}-a}{2^{-p-1}-(B+1-A_i)2^{-q}} \leq \frac{2^{-p-1}+i2^{-p}-a}{A_i2^{-q}}$  because  $2^{-p-1} - (B+1)2^{-q} \geq 0$  since  $(B+1)2^{-q} \leq 2^\omega 2^{-q} \leq 2^{-p-1}$ . Consequently, for the first coordinate we get that  $f_{\mathcal{I}}(a, b)_1 \leq a + 2^{-p} + A_i 2^{-q} \frac{2^{-p-1}+i2^{-p}-a}{A_i 2^{-q}} \leq (i+1)2^{-p} + 2^{-p-1}$ . Also, since  $i2^{-p} \leq a \leq i2^{-p} + 2^{-p-1}$ , it is clear that  $f_{\mathcal{I}}(a, b)_1 \geq (i+1)2^{-p}$ . So finally,  $f_{\mathcal{I}}(R_{i,3}^{lin}) \subseteq [(i+1)2^{-p}, (i+1)2^{-p}+2^{-p-1}] \times \{0\} \subset R_{i+1}$ .
- If  $(a, b) \in R_{i,3}^{sat}$ : the image of the second component is always 0 so it's easy for this one, also from Definition 23,  $b\beta^i - 3 \geq \frac{2^{-p-1}+i2^{-p}-a}{2^{-p-1}-(B+1-A_i)2^{-q}} \geq \frac{2^{-p-1}+i2^{-p}-a}{2^{-p-1}-(B+1)2^{-q}}$  because  $A_i \geq 0$ . Consequently, for the first coordinate we get that  $f_{\mathcal{I}}(a, b)_1 \leq (i+1)2^{-p} + 2^{-p-1} - (2^{-p-1} - (B+1)2^{-q}) \frac{2^{-p-1}+i2^{-p}-a}{2^{-p-1}-(B+1)2^{-q}} \leq (i+1)2^{-p} + 2^{-p-1} + i2^{-p} + 2^{-p-1} - a \leq (i+1)2^{-p} + 2^{-p-1}$  since  $a \leq i2^{-p} + 2^{-p-1}$ . Also since  $b\beta^i - 3 \leq 1$  we get that  $f_{\mathcal{I}}(a, b)_1 \geq (i+1)2^{-p} + 2^{-p-1} - (2^{-p-1} - (B+1)2^{-q}) \times 1 \geq (i+1)2^{-p} + (B+1)2^{-q}$ . So finally,  $f_{\mathcal{I}}(R_{i,3}^{sat}) \subseteq [(i+1)2^{-p} + (B+1)2^{-q}, (i+1)2^{-p} + 2^{-p-1}] \times \{0\} \subset R_{i+1}$ .
- If  $(a, b) \in R_{i,4}^{lin} = R_{i,1}^{lin}$ : then  $f_{\mathcal{I}}(a, b) = (a + 2^{-p} + A_i 2^{-q}, b - 1^* \beta^{-i})$  so  $f_{\mathcal{I}}(R_{i,4}^{lin}) = f_{\mathcal{I}}([i2^{-p}, i2^{-p} + (B+1 - A_i)2^{-q}] \times [4\beta^{-i}, 5\beta^{-i}]) = [(i+1)2^{-p} + A_i 2^{-q}, (i+1)2^{-p} + (B+1)2^{-q}] \times [0, \beta^{-i}] \subseteq R_{i+1}$  because  $(B+1)2^{-q} \leq 2^{-p-1}$ .
- If  $(a, b) \in R_{i,4}^{sat}$ : then  $f_{\mathcal{I}}(a, b) = ((i+1)2^{-p} + (B+1)2^{-q}, 0)$  so  $f_{\mathcal{I}}(R_{i,4}^{sat}) = \{(i+1)2^{-p} + (B+1)2^{-q}\} \times \{0\} \subseteq R_{i+1}$ .

□

We now get to the core lemma of the simulation. Up to this point, we were only interested in forward simulation: that is given a point, what are the iterates of  $x$ . In order to prove the NP-hardness result, we need a backward result: given a point, what are the possible preimages of it. To this end, we introduce new subregions  $R_i^{unsat}$  of the  $R_i$ , that we call *unsaturated*. Intuitively,  $R_i^{unsat}$  corresponds to the encodings where  $\sigma \leq B$ , that is the sum did not saturate at  $B+1$ . We also introduce the  $R_{fin}$  region, that will be the region to reach. We will be interested in the preimages of  $R_{fin}$ .

**Definition 27 (Unsaturated regions).** For  $i \in \{1, \dots, n+1\}$ , define

$$R_i^{unsat} = [i2^{-p}, i2^{-p} + B2^{-q}] \times [\beta^{-n-1}, \beta^{-i+1} - \beta^{-n-1}]$$

$$R_{fin} = [(n+1)2^{-p} + B2^{-q} - 2^{-q-1}, (n+1)2^{-p} + B2^{-q}] \times [\beta^{-n-1}, 2\beta^{-n-1}]$$

**Lemma 28 (Simulation is reversible).** Let  $i \in \{2, \dots, n\}$  and  $(a, b) \in R_i^{unsat}$ . Then the only points  $\mathbf{x}$  such that  $f_{\mathcal{I}}(\mathbf{x}) = (a, b)$  are:

- $\mathbf{x} = (a - 2^{-p}, b + 0^* \beta^{-i+1}) \in R_{i-1,0^*} \cap R_{i-1}^{unsat}$
- $\mathbf{x} = (a - 2^{-p}, \beta^i - b + 0^* \beta^{-i+1}) \in R_{i-1,2} \cap R_{i-1}^{unsat}$



- $\mathbf{x} = (a - 2^{-p} - A_i 2^{-q}, b + 1^* \beta^{-i+1}) \in R_{i-1,1^*}^{lin} \cap R_{i-1}^{unsat}$  (only if  $a \geq 2^{-p} + A_i 2^{-q}$ )

*Proof.* First notice that since  $f_{\mathcal{I}}(R_i) \subseteq R_{i+1}$  for all  $i \in \{0, \dots, n\}$ , the only candidates for  $\mathbf{x}$  must belong to  $R_{i-1}$ . Furthermore, on each affine region, there can only be one candidate except if the function is trivial.

A close look at the proof of Lemma 26 reveals that:

- $f_{\mathcal{I}}(R_{i-1,0}) \subseteq [(i+1)2^{-p}, (i+1)2^{-p} + 2^{-p-1}] \times \{0\}$  shareing no point with  $R_i^{unsat}$  so there is no possible candidate
- $f_{\mathcal{I}}(R_{i-1,1}) = R_i$  and there is only one possible candidate
- $f_{\mathcal{I}}(R_{i-1,2}) = R_i$  and there is only one possible candidate
- $f_{\mathcal{I}}(R_{i-1,3}) \subseteq [(i+1)2^{-p}, (i+1)2^{-p} + 2^{-p-1}] \times \{0\}$  so like  $R_{i-1,0}$  there is no possible candidate
- $f_{\mathcal{I}}(R_{i-1,4}^{lin}) \subseteq R_i$  and there is only one possible candidate
- $f_{\mathcal{I}}(R_{i-1,4}^{sat}) \subseteq [(i+1)2^{-p} + (B+1)2^{-q}, (i+1)2^{-p} + 2^{-p-1}] \times [0, \beta^{-i}]$  sharing no point with  $R_i^{unsat}$  so there is no possible candidate

It is then only a matter of checking that the claimed formulas work and they trivially do except for the case of  $R_{i-1,4}^{lin}$  where we need the potential candidate to belong to the region.  $\square$

The goal of those results is to show that if there is a point in  $R_{fin}$  that is reachable from  $R_0$  then we can extract, from its trajectory, a configuration that also reaches  $R_{fin}$ . Furthermore, we arranged so that  $R_{fin}$  contains the encoding of only one configuration:  $(n+1, B)$  (see Lemma 18).

**Lemma 29 (Backward-forward identity).** *For any point  $\mathbf{x} \in R_{fin}$ , if there exists a point  $\mathbf{y} \in R_0$  and an integer  $k$  such that  $f_{\mathcal{I}}^{[k]}(\mathbf{y}) = \mathbf{x}$  then there exists a configuration  $c = (1, 0, \varepsilon_1, \dots, \varepsilon_n)$  such that  $f_{\mathcal{I}}^{[k]}(c) \in R_{fin}$ .*

*Proof.* Define  $\mathbf{y}_0 = \mathbf{y}$  and  $\mathbf{y}_{i+1} = f_{\mathcal{I}}(\mathbf{y}_i)$  for all  $i \in \{0, k-1\}$ . Since  $\mathbf{y}_0 \in R_0$ , we immediately get that  $\mathbf{y}_i \in R_i$  using Lemma 26 and in particular,  $k \geq n+1$  because  $\mathbf{y}_k = \mathbf{x} \in R_{n+1,0^*}$ .

Now apply Lemma 28 starting from  $\mathbf{y}_{n+1} \in R_{n+1}^{unsat}$ : we conclude that for all  $i \geq 1$ ,  $\mathbf{y}_i \in (R_{i,1} \cup R_{i,2} \cup R_{i,4}^{lin}) \cap R_n^{unsat}$ . Define  $\varepsilon_i = 0$  if  $\mathbf{y}_i \in R_{i,1} \cup R_{i,2}$  and 1 if  $\mathbf{y}_i \in R_{i,4}^{lin}$ . Write  $\mathbf{y}_i = (a_i, b_i)$ . Again using Lemma 26 we get that  $a_{i-1} = a_i - 2^{-p} - \varepsilon_i A_i 2^{-q}$  (just check all three cases). Also since  $\mathbf{x} = \mathbf{y}_{n+1} \in R_{fin}$  then  $a_{n+1} \in [(n+1)2^{-p} + B2^{-q} - 2^{-q-1}, (n+1)2^{-p} + B2^{-q}]$ . Finally,  $\mathbf{y}_0 \in R_0$  so  $f_{\mathcal{I}}(a_0, b_0) = (2^{-p}, 0) = (a_1, b_1)$ . We conclude that  $a_1 = 2^{-p}$ . Putting everything together we get:

$$\begin{cases} a_{n+1} = (n+1)2^{-p} + 2^{-q} \sum_{i=1}^n \varepsilon_i A_i \\ a_{n+1} \in [(n+1)2^{-p} + B2^{-q} - 2^{-q-1}, (n+1)2^{-p} + B2^{-q}] \end{cases}$$

Since the  $A_i, B$  are integers and  $\varepsilon_i \in \{0, 1\}$ , we get that  $B = \sum_{i=1}^n A_i \varepsilon_i$ . Apply Lemma 22 on the configuration to conclude.  $\square$

**Lemma 30 (Final region is accepting).** *For any configuration  $c$ , if  $\langle c \rangle \in R_{fin}$  then  $c = (n + 1, B)$ .*

*Proof.* Write  $c = (i, \sigma, \varepsilon_i, \dots, \varepsilon_n)$ , then  $\langle c \rangle = \left( i2^{-p} + \sigma2^{-q}, \sum_{j=i}^n \varepsilon_j^* \beta^{-j} + 0^* \beta^{-n-1} \right)$ . It implies that  $i2^{-p} + \sigma2^{-q} \in [(n + 1)2^{-p} + B2^{-q} - 2^{-q-1}, (n + 1)2^{-p} + B2^{-q}]$  and because  $i$  is an integer in range  $[0, n + 1]$  and  $\sigma$  an integer in range  $[0, B + 1]$ , necessarily  $i = n + 1$  and  $\sigma = B$ .  $\square$

### 3.4 Complexity result

We now have all the tools to show that REACH-REGION-TIME is an NP-hard problem.

**Theorem 31.** *REACH-REGION-TIME is NP-hard for  $d \geq 2$ .*

*Proof.* Let  $\mathcal{I} = (B, A_1, \dots, A_n)$  be an instance of SUBSET-SUM. We consider the instance  $\mathcal{J}$  of REACH-REGION-TIME defined in the previous section with maximum number of iterations set to  $n$  (the number of  $A_i$ ), the initial region set to  $R_0$  and the final region set to  $R_{fin}$ . One easily checks that this instance has polynomial size in the size of  $\mathcal{I}$ . The two directions of the proof are:

- If  $\mathcal{I}$  is satisfiable then use Lemma 17 and Lemma 22 to conclude that there is a point  $x \in R_0$  in the initial region such that  $f_{\mathcal{I}}^{[n]}(x) \in R_{fin}$  so  $\mathcal{J}$  is satisfiable.
- If  $\mathcal{J}$  is satisfiable then there exists  $x \in R_0$  and  $k \leq n$  such that  $f_{\mathcal{I}}^{[k]}(x) \in R_{fin}$ . Use Lemma 29 and Lemma 22 to conclude that there exists a configuration  $c = (1, 0, \varepsilon_1, \dots, \varepsilon_n)$  such that  $\langle T_{\mathcal{I}}^{[k]}(c) \rangle = f_{\mathcal{I}}^{[k]}(\langle c \rangle) \in R_{fin}$ . Apply Lemma 30 and use the injectivity of the encoding to conclude that  $T_{\mathcal{I}}^{[k]}(c) = (n + 1, B)$  and Lemma 18 to get that  $\mathcal{I}$  is satisfiable.  $\square$

## 4 Bounded Time Reachability is in NP

In the previous section we have shown that the REACH-REGION-TIME problem is NP-hard. We now give a more precise characterization of the complexity of this problem, by proving that it is NP-complete. Since we have shown its NP-hardness, the only thing that remains to be shown is that REACH-REGION-TIME belongs to NP. This is done in this section.

### 4.1 Notations and definitions

For any  $i \in \{1, \dots, d\}$ , let  $\pi_i^d: \mathbb{R}^d \rightarrow \mathbb{R}$  denote the  $i^{th}$  projection function, that is,  $\pi(x_1, \dots, x_d) = x_i$ . Let  $g_d: \mathbb{R}^{d+1} \rightarrow \mathbb{R}^d$  be defined by  $g_d(x_1, \dots, x_{d+1}) = (x_1, \dots, x_d)$ . For a square matrix  $A$  of size  $(d + 1) \times (d + 1)$  define the following pair of projection functions. The first function  $h_{1,d}$  takes as input a square matrix

$A$  of size  $(d+1) \times (d+1)$  and returns a square matrix of size  $d \times d$  that is the upper-left block of  $A$ . The second function  $h_{2,d}$  takes as input a square matrix  $A$  of size  $(d+1) \times (d+1)$  and returns the vector of size  $d$  given by  $[a_{1,d+1} \cdots a_{d,d+1}]^T$  (the last column of  $A$  minus the last element).

Let  $s$  denote the size function, its domain of objects will be overloaded and understood from the context. For  $x \in \mathbb{Z}$ ,  $s(x)$  is the length of the encoding of  $x$  in base 2. For  $x \in \mathbb{Q}$  with  $x = \frac{p}{q}$  with  $p$  and  $q$  coprime, we have  $s(x) = \max(s(p), s(q))$ . For an affine function  $f$  we define the size of  $f(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$  (where all entries of  $A$  and  $\mathbf{b}$  are rationals) as:  $s(f) = \max(\max_{i,j}(s(a_{i,j})), \max(s(b_i)))$ . We define the size of a polyhedron  $r$  defined by  $A\mathbf{x} \leq \mathbf{b}$  as:  $s(r) = \max(s(A), s(\mathbf{b}))$ .

We define the size of a piecewise affine function  $f$  as:  $s(f) = \max_i(s(f_i), s(r_i))$  where  $f_i$  denotes the restriction of  $f$  to  $r_i$  the  $i^{\text{th}}$  region.

We define the *signature* of a point  $\mathbf{x}$  as the sequence of indices of the regions traversed by the iterates of  $f$  on  $\mathbf{x}$  (that is, the region trajectory).

## 4.2 REACH-REGION-TIME is in NP

In order to solve a reachability problem, we will formulate it with linear algebra. However a crucial issue here is that of the size of the numbers, especially when computing powers of matrices. Indeed, if taking the  $n^{\text{th}}$  power of  $A$  yields a representation of exponential size, no matter how fast our algorithm is, it will run on exponentially large instances and thus be slow.

First off, we show how to move to homogeneous coordinates so that  $f$  becomes piecewise linear instead of piecewise affine.

**Lemma 32.** *Assume that  $f(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$  with  $A = (a_{i,j})_{1 \leq i,j \leq d}$  and let  $y = A'(\mathbf{x}, 1)^T$  where  $A'$  is the block matrix  $\begin{pmatrix} A & \mathbf{b} \\ 0 & 1 \end{pmatrix}$ . Then  $f(x) = g_d(A'(\mathbf{x}, 1)^T)$ .*

*Remark 33.* Notice that this lemma extends nicely to the composition of affine functions: if  $f(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$  and  $h(\mathbf{x}) = C\mathbf{x} + \mathbf{d}$  then  $h(f(x)) = g_d(C'A'(\mathbf{x}, 1)^T)$ .

We can now state the main lemma, namely that the size of the iterates of  $f$  vary linearly in the number of iterates, assuming that  $f$  is piecewise affine.

**Lemma 34.** *Let  $d \geq 2$  and  $f \in PAF_d$ . Assume that all the coefficients of  $f$  on all regions are rationals. Then for all  $t \in \mathbb{N}$ ,  $s(f^{[t]}) \leq (d+1)^2 s(f)pt + (t-1)[\log_2(d+1)]$  where  $p$  is the number of regions of  $f$ . This inequality holds even if all rationals are taken to have the same denominator.*

*Proof.* Using Lemma 32, we get that  $f^{[t]}(\mathbf{x}) = g_d(h^{[t]}([\mathbf{x} \ 1]^T))$ , where  $h$  is a piecewise linear function in dimension  $d+1$  such that  $s(h) = s(f)$ . We show this result by induction on  $t$  for  $h$ . The result then follows for  $f$ . In all cases we take all rationals to have the same denominator.

In the case  $t = 1$ , it suffices to see that taking all rationals to have the same denominator involves multiplying the numerator and denominators by at most the lowest common multiple of all numbers, and hence is at most  $2^{s(f)(p(d+1)^2)}$ .

Indeed the greatest number is  $2^{s(h)}$  by definition, and there are  $(d+1)^2$  numbers per region (the entries of the matrix).

Assume the result is true for  $t \in \mathbb{N}$ . Let  $\mathbf{y} \in \mathbb{Q}^{d+1}$ . Then  $h^{[t+1]}(\mathbf{y}) = B_{t+1} \cdots B_1 \mathbf{y}$ , where  $B_i$ 's are the matrices corresponding to some regions of  $h$ . In particular,  $s(B_i) \leq s(h)$ . From the induction hypothesis we can assume that all rationals have the same denominator and we get that  $s(B_t \cdots B_1) \leq (d+1)^2 s(h) p t + (t-1) \lceil \log_2(d+1) \rceil$ . It follows<sup>6</sup> that for any  $1 \leq i, j \leq d+1$ :

$$\begin{aligned} s((B_{t+1} \cdots B_1)_{i,j}) &= s\left(\sum_{k=1}^{d+1} (B_{t+1})_{i,k} (B_t \cdots B_1)_{k,j}\right) \\ &\leq \lceil \log_2(d+1) \rceil + s(B_{t+1}) + s(B_t \cdots B_1) \\ &\leq \lceil \log_2(d+1) \rceil + s(h) + (d+1)^2 s(h) p t + (t-1) \lceil \log_2(d+1) \rceil \\ &\leq (d+1)^2 s(h) p (t+1) + t \lceil \log_2(d+1) \rceil \end{aligned}$$

This shows the result for the particular region where  $y$  belongs. Since the bound does not depend on  $y$  and  $h^{[t+1]}$  has finitely many regions, it is true for all regions of  $h^{[t+1]}$ .  $\square$

Finally, we need some result about the size of solutions to systems of linear inequalities. Indeed, if we are going to quantify over the existence of a solution of polynomial size, we must ensure that the size constraints do not change the satisfiability of the system.

**Lemma 35 ([11]).** *Let  $A$  be a  $N \times d$  integer matrix and  $\mathbf{b}$  an integer vector. If the  $A\mathbf{x} \leq \mathbf{b}$  system admits a solution, then there exists a rational solution  $x_s$  such that  $s(x_s) \leq (d+1)L + (2d+1) \log_2(2d+1)$  where  $L = \max(s(A), s(b))$ .*

*Proof.* See Theorem 5 of [11]:  $s(x_s) \leq s((2d+1)! 2^{L(2d+1)})$ .  $\square$

Putting everything together, we obtain a fast nondeterministic algorithm to solve REACH-REGION-TIME. The nondeterminism allows us to choose a signature for the solution. Once the signature is fixed, we can write it as a linear program of reasonable size using Lemma 34 and solve it. The remaining issue is the one of the size of solution but fortunately Lemma 35 ensures us that there is a small solution that can be found quickly.

**Theorem 36.** *REACH-REGION-TIME is in NP.*

*Proof.* The idea of the proof is to nondeterministically choose a signature for a solution, that is a sequence of regions for the iterates of the solution. We then build a system of linear inequalities stating that a point  $\mathbf{x}$  belongs to the initial region and that the iterates match the signature chosen and finally that the iterates reach the final region. Using the results of the previous section, we can build this system in polynomial time and solve it in non-deterministic polynomial time. Here is an outline of the algorithm:

<sup>6</sup> Use elementary properties of the size function:  $s(xy) \leq s(x) + s(y)$ ,  $s(x_1 + \cdots + x_k) \leq s(k) + \max_k s(x_k)$

- Non-deterministically choose  $t \leq T$
- Non-deterministically choose regions  $r_1, \dots, r_{t-1}$  regions of  $f$
- Define  $r_0 = R_0$  the initial region and  $r_t = R$  the final region
- Build  $(S)$  the system  $A\mathbf{x} \leq \mathbf{b}$  stating that the signature of  $\mathbf{x}$  matches  $r$
- Non-deterministically choose  $\mathbf{x}_s$  a rational of polynomial size in the size of  $(S)$
- Accept if  $A\mathbf{x}_s \leq \mathbf{b}$

We have two things to prove. First we need to show that this algorithm indeed has non-deterministic polynomial running time. Second we need to show that it is correct. Recall that  $T$  is a unary input of the problem.

The complexity of the algorithm is clear, assuming that  $(S)$  is of polynomial size. Indeed verifying that a rational point satisfies a system of linear inequalities with rational coefficients can be done in polynomial time.

We build  $(S)$  this way:  $(S) = \cup_{i=1}^t (S_i)$  where  $(S_i)$  states that  $f^{[i]}(\mathbf{x}) \in r_i$ . Since we choose a signature of  $\mathbf{x}$  we know that if  $\mathbf{x}$  satisfies the system then from Lemma 32  $f^{[i]}(\mathbf{x}) = g_d(A'_{i-1} \cdots A'_1(\mathbf{x}, 1)^T)$  where  $A'_j$  is the matrix corresponding to the region  $r_j$ . Write  $C_i = A'_{i-1} \cdots A'_1$  and define  $(S_i)$  by the system  $g_d(C_i(\mathbf{x}, 1)^T) \in r_i$ . Since  $r_i$  is a polyhedron,  $(S_i)$  is indeed a system of linear inequalities<sup>7</sup>.

We can now see that  $S$  is of polynomial size using Lemma 34. Indeed,  $s(C_i) \leq s(f^{[i]}) \leq \text{poly}(s(f), i)$ , thus  $s((S_i)) \leq s(C_i) + s(r_i) \leq \text{poly}(s(f), i)$  because the description of the regions is part of the size of  $f$ . And finally  $s((S)) \leq \text{poly}(s(f), t)$ .

The correctness follows from the construction of the system and Lemma 35. More precisely we show that  $\mathbf{x} \in (S)$  if and only if  $\forall i \in \{0, \dots, t\}, f^{[i]}(\mathbf{x}) \in r_i$ . Indeed,  $(S) \Leftrightarrow \forall i \in \{0, \dots, t\}, \mathbf{x} \in (S_i)$  and by definition  $(S_i) \Leftrightarrow f^{[i]}(\mathbf{x}) \in r_i$  since  $g_d(C_i(\mathbf{x}, 1)^T) = f^{[i]}(\mathbf{x})$ . Then by Lemma 35, we get that  $\exists \mathbf{x} \in (S) \Leftrightarrow \exists \mathbf{x} \in (S)$  and  $s(\mathbf{x}) \leq \text{poly}(s((S)))$ .  $\square$

## 5 Other Bounded Time Results

In this section, we give succinct proofs of the other result mentioned in the introduction about CONTROL-REGION-TIME. The proof is based on the same arguments as before.

**Theorem 37.** *Problem CONTROL-REGION-TIME is coNP-hard for  $d \geq 2$ .*

*Proof.* The proof is exactly the same except for two details:

- we modify  $f$  over  $R_{n+1}$  as follows: divide  $R_{n+1}$  in three regions:  $R_{low}$  that is below  $R_{fin}$ ,  $R_{fin}$  and  $R_{high}$  that is above  $R_{fin}$ . Then build  $f$  such that  $f(R_{low}) \subseteq R_{low}$ ,  $f(R_{fin}) \subseteq R_{fin}$  and  $f(R_{high}) \subseteq R_{low}$ .
- we choose a new final region  $R'_{fin} = R_{low}$ .

<sup>7</sup> More precisely if  $r_i$  is defined by  $P_i(\mathbf{x}, 1)^T \leq 0$  then  $(S_i)$  is the system  $P_i C_i(\mathbf{x}, 1)^T \leq 0$

Let  $\mathcal{I} = (B, A_1, \dots, A_n)$  be an instance of NOSUBSET-SUM, let  $\mathcal{J}$  be the corresponding instance of CONTROL-REGION-TIME we just built. We have to show that  $\mathcal{I}$  has no subset sum if and only if  $\mathcal{J}$  is “controlled”. This is the same as showing that  $\mathcal{I}$  has a subset sum if and only if  $\mathcal{J}$  has points never reaching  $R'_{fin}$ .

Now assume for a moment that the instance is in SUBSET-SUM (as opposed to NOSUBSET-SUM), then by the same reasoning as the previous proof, there will be a point that reaches the old  $R_{fin}$  region (and is disjoint from  $R'_{fin}$ ). And since  $R_{fin}$  is a  $f$ -stable region, this point will never reach  $R'_{fin}$ .

And conversely, if the control problem is not satisfied, necessarily there is a point whose trajectory went through the old  $R_{fin}$  (otherwise it would have reached either  $R_{low} = R'_{fin}$  or  $R_{high}$  but  $f(R_{high}) \subseteq R_{low}$ ). Now we proceed as in the proof of Theorem 31 to conclude that there is a subset that sums to  $B$ , and thus  $\mathcal{I}$  is satisfiable.  $\square$

**Theorem 38.** *Problem CONTROL-REGION-TIME is in coNP for  $d \geq 2$ .*

*Proof.* Again the proof is very similar to that of Theorem 36: we have to build a non-deterministic machine that accepts the “no” instances. The algorithm is exactly the same except that we only choose signatures that avoid the final region (as opposed to ending in the final region) and are of maximum length (that is  $t = T$  as opposed to  $t \leq T$ ). Indeed, if there is such a trajectory, the problem is not satisfied. And for the same reasons as Theorem 36, it runs in non-deterministic polynomial time.  $\square$

## 6 Fixed Precision Results

**Theorem 39.** *REACH-REGION-PRECISION and CONTROL-REGION-PRECISION are PSPACE-hard.*

*Proof.* Consider a polynomial space Turing machine  $\mathcal{M} = (Q, \Sigma, B, q_0, F, \delta)$ . Without loss of generality, we can assume that  $F = \{q_f\} \subset Q$  (there is a single accepting state) and that the working alphabet is  $\Sigma = \{0, 1, 2, \dots, \beta\}$ , assuming that  $B = 0$  is the blank character, and  $\delta : Q \times \Sigma \rightarrow Q \times \Sigma \times \{\triangleleft, \square, \triangleright\}$  is complete. We also assume the set of internal states to be such that  $Q \subseteq \Sigma^m$  for some  $m$ .

Let  $c$  be an instantaneous configuration (sometimes also called *ID* for *Instantaneous Description*) of  $\mathcal{M}$ . We write  $c = (x, \sigma, q, y)$  where  $x$  (resp.  $y$ ) is encoding the part of the tape on the left (resp. right) of the head,  $\sigma$  is the symbol under the head, and  $q$  is the state of the machine. Specifically, if the non-blank part of the tape is  $s_{-n} \dots s_{-1} \sigma s_1 s_2 \dots s_m$ , with the head in front of  $\sigma$ , then  $x$  is encoded as word  $s_{-1} s_{-2} \dots s_{-n}$ , and  $y$  as word  $s_1 s_2 \dots s_m$ .

Define the encoding of configuration  $c$  as  $\langle c \rangle = (\langle x \rangle, \langle q\sigma y \rangle)$  where for any word  $w \in \Sigma^*$ ,  $\langle w \rangle = \sum_{i=1}^{|w|} 2w_i(2^\gamma)^{-i}$ , where  $\gamma$  is such that  $2\beta + 1 \leq 2^\gamma$ . Define regions  $R_{\alpha, q, \sigma} = [\langle \alpha \rangle] \times [\langle q\sigma \rangle]$  where  $[\langle w \rangle]$  is a shortcut for  $[\langle w \rangle] = [\langle w \rangle + (2^\gamma)^{-|w|}]$ . Intuitively,  $R_{\alpha, q, \sigma}$  contains all configurations in state  $q$ , with symbol  $\sigma$  under the head and symbol  $\alpha$  immediately at the left of the head. By construction

$\langle c \rangle \in R_{x_1, q, \sigma}$  with the above notations. Finally, for  $\alpha, \sigma \in \Sigma$ ,  $q \in Q$ , define  $f$  on region  $R_{\alpha, q, \sigma}$  by:

$$f(a, b) = \begin{cases} a2^{-\gamma} + \langle \sigma' \rangle, (b - \langle q\sigma \rangle)2^\gamma + \langle q' \rangle & \text{if } \delta(q, \sigma) = (q', \sigma', \triangleright) \\ (a, b - \langle q\sigma \rangle + \langle q'\sigma' \rangle) & \text{if } \delta(q, \sigma) = (q', \sigma', \square) \\ (a - \langle \alpha \rangle)2^\gamma, (b - \langle q\sigma \rangle)2^{-\gamma} + \langle q'\alpha\sigma' \rangle & \text{if } \delta(q, \sigma) = (q', \sigma', \triangleleft) \end{cases}$$

It is clear from the definition that  $f$  is piecewise affine over its domain of definition. Let  $T$  be the function corresponding to one step of computation of  $\mathcal{M}$ :  $T$  is acting on configurations and maps any configuration  $c$  to the corresponding next configuration according to the program of  $\mathcal{M}$ . A simple case analysis shows that for any configuration  $c$ ,  $\langle T(c) \rangle = f(\langle c \rangle)$ , using the fact that the blank character  $B$  was chosen to be 0.

Observe furthermore that by the choice of the encoding, all the regions  $R_{\alpha, q, \sigma}$  are closed and at positive distance from each other:  $R_{\alpha, q, \sigma} \cap R_{\alpha', q', \sigma'} = \emptyset$  whenever  $(\alpha, q, \sigma) \neq (\alpha', q', \sigma')$ . It follows that  $f$  can be easily extended to a continuous piecewise linear function defined over the whole domain  $[0, 1]$  (similar arguments are used in [12]). By construction it will still satisfy that for any configuration  $c$ ,  $\langle T(c) \rangle = f(\langle c \rangle)$ .

We can now state the reduction from problem **LINSPACE-WORD**: consider an instance  $(\mathcal{M}, w)$  of this decision problem. Define  $\varepsilon = (2^\gamma)^{-(|w|+m+1)}$  where the choice of  $m$  was explained above. Then for any configuration  $c$  reachable from the initial configuration  $c_0 = (\varepsilon, q_0, w_1, w_2 \cdots w_{|w|})$ , we have the stronger property that  $\langle c \rangle = \lfloor \frac{\langle c_0 \rangle}{\varepsilon} \rfloor \varepsilon$ . Indeed, by assumption the machine never uses more space than the size of the input, thus the left and right part of tape are always smaller than  $|w|$  at any point during the computation, and we simply need an extra  $m+1$  space to store the current state of the machine. In other words, rounding to  $\varepsilon$  does not perturbate the simulation. Consequently, we get that for any reachable configuration  $c$ ,  $\langle T(c) \rangle = f_\varepsilon(\langle c \rangle)$ .

Define  $R_0 = \{\langle c_0 \rangle\}$  and  $R = [0, 1] \times [\langle q_f \rangle]$  that are convex regions. Then the instance  $(f, R_0, R, \varepsilon)$  of **REACH-REGION-PRECISION** is satisfiable if and only if  $(\mathcal{M}, w)$  belongs to problem **LINSPACE-WORD**. One easily checks that  $(f, R_0, R, \varepsilon)$  has polynomial size in the size of  $(\mathcal{M}, w)$ .

The same instance also works for **CONTROL-REGION-PRECISION**. If we want to make  $R_0$  a region with non-empty interior, just take a ball of radius smaller than  $\varepsilon(2^\gamma)^{-2}$  around  $\langle c_0 \rangle$  so that any input error is removed after the first application of the function  $f_\varepsilon$ .  $\square$

**Theorem 40.** *REACH-REGION-PRECISION and CONTROL-REGION-PRECISION are in PSPACE.*

*Proof.* Let  $N = \lfloor \varepsilon^{-1} \rfloor$  and consider the graph  $G = (V, E)$  where

$$V = \{0, \dots, N-1\}^d \quad C_\alpha = \prod_{k=1}^d [\alpha_k \varepsilon, (\alpha_k + 1)\varepsilon[ \quad (\alpha \in V)$$

$$E = \{(\alpha, \beta) \mid f(C_\alpha) \cap C_\beta \neq \emptyset\} \quad S = \{\alpha \mid f(R_0) \cap C_\alpha \neq \emptyset\}$$

$$T = \{\alpha \mid R \cap C_\alpha \neq \emptyset\}$$

We can now restate our reachability problem in the graph  $G$  as an accessibility problem from  $S$  to  $T$ . This can be done in space logarithmic in the size of the graph, using the fact that accessibility in a graph with  $N$  vertices can be done in non-deterministic space  $O(\log N)$ , and using the fact that  $NSPACE(N) = SPACE(N^2)$  (Savitch's Theorem) [16, Theorem 8.5]. Since the graph is of size  $O(N^d)$ , this requires space  $O(\log^2 N) = O(-\log^2 \varepsilon) = O(n^2)$  if  $\varepsilon = 2^{-n}$ . Also note that computing the transitions of the graph is fast since  $f$  is a piecewise affine function. The same proof applies to CONTROL-REGION-PRECISION except we now want to know if for every vertex in  $S$  there is a path to  $T$ .  $\square$

## References

1. Asarin, E., Maler, O., Pnueli, A.: Reachability analysis of dynamical systems having piecewise-constant derivatives. *Theoretical Computer Science* 138(1), 35–65 (Feb 1995)
2. Asarin, E., Schneider, G.: Widening the boundary between decidable and undecidable hybrid systems. In: Brim, L., Jancar, P., Kretínský, M., Kucera, A. (eds.) CONCUR 2002 - Concurrency Theory, 13th International Conference, Brno, Czech Republic, August 20-23, 2002, Proceedings. *Lecture Notes in Computer Science*, vol. 2421, pp. 193–208. Springer (2002), <http://link.springer.de/link/service/series/0558/bibs/2421/24210193.htm>
3. Asarin, E., Schneider, G., Yovine, S.: On the decidability of the reachability problem for planar differential inclusions. In: Benedetto, M.D.D., Sangiovanni-Vincentelli, A.L. (eds.) Hybrid Systems: Computation and Control, 4th International Workshop, HSCC 2001, Rome, Italy, March 28-30, 2001, Proceedings. *Lecture Notes in Computer Science*, vol. 2034, pp. 89–104. Springer (2001), <http://link.springer.de/link/service/series/0558/bibs/2034/20340089.htm>
4. Bazille, H., Bournez, O., Gomaa, W., Pouly, A.: On the complexity of bounded time reachability for piecewise affine systems. In: Ouaknine, J., Potapov, I., Worrell, J. (eds.) Reachability Problems, *Lecture Notes in Computer Science*, vol. 8762, pp. 20–31. Springer International Publishing (2014), [http://dx.doi.org/10.1007/978-3-319-11439-2\\_2](http://dx.doi.org/10.1007/978-3-319-11439-2_2)
5. Bell, P., Chen, S.: Reachability problems for hierarchical piecewise constant derivative systems. In: Abdulla, P., Potapov, I. (eds.) Reachability Problems, *Lecture Notes in Computer Science*, vol. 8169, pp. 46–58. Springer Berlin Heidelberg (2013), [http://dx.doi.org/10.1007/978-3-642-41036-9\\_6](http://dx.doi.org/10.1007/978-3-642-41036-9_6)
6. Ben-Amram, A.M.: Mortality of iterated piecewise affine functions over the integers: Decidability and complexity. In: STACS. pp. 514–525 (2013)
7. Blondel, V.D., Bournez, O., Koiran, P., Tsitsiklis, J.: The stability of saturated linear dynamical systems is undecidable. *Journal of Computer and System Science* 62(3), 442–462 (May 2001), <http://dx.doi.org/10.1006/jcss.2000.1737>
8. Garey, M.R., Johnson, D.S.: *Computers and Intractability*. W. H. Freeman and Co (1979)
9. Henzinger, T.A., Kopke, P.W., Puri, A., Varaiya, P.: What's decidable about hybrid automata? *Journal of Computer and System Sciences* 57(1), 94–124 (Aug 1998)



10. Karp, R.M.: Reducibility among combinatorial problems. Springer (1972)
11. Koiran, P.: Computing over the reals with addition and order. *Theor. Comput. Sci.* 133(1), 35–47 (1994)
12. Koiran, P., Cosnard, M., Garzon, M.: Computability with low-dimensional dynamical systems. *Theoretical Computer Science* 132(1-2), 113–128 (Sep 1994)
13. Koiran, P., Cosnard, M., Garzon, M.: Computability with Low-Dimensional Dynamical Systems. *Theoretical Computer Science* 132, 113–128 (1994)
14. Moore, C.: Generalized shifts: unpredictability and undecidability in dynamical systems. *Nonlinearity* 4(3), 199–230 (1991)
15. Siegelmann, H.T., Sontag, E.D.: On the computational power of neural nets. *Journal of Computer and System Sciences* 50(1), 132–150 (Feb 1995)
16. Sipser, M.: *Introduction to the Theory of Computation*. PWS Publishing Company (1997)