

Positive Invariance of Constrained Affine Dynamics and Its Applications to Hybrid Systems and Safety Verification

Jinglai Shen*

Original: May 22, 2008; Revised: December 15, 2008

Abstract

Motivated by long-time dynamic analysis of hybrid systems and safety verification problems, this paper addresses fundamental positive invariance issues of an affine dynamics on a general convex polyhedron and their applications. Necessary and sufficient algebraic conditions are established for the existence of a positively invariant set of an affine system on a polyhedron using the recent tools of lexicographic relation and long-time oscillatory dynamic analysis. It is shown that these existence conditions can be finitely verified once the eigenvalues of the defining matrix are known. The positive invariance results are applied to obtain an explicit characterization of global switching behaviors of piecewise affine systems (PASs). In particular, the PASs with isolated equilibria and infinite mode switchings are characterized via the positive invariance conditions. The positive invariance techniques developed in this paper are also exploited to show decidability of safety verification of affine dynamics on semialgebraic sets.

1 Introduction

Modeling, analysis, control, and computation of hybrid dynamical systems has received fast growing interest in the past decade, driven by various important applications such as air traffic control [29], nonsmooth mechanical/electrical systems [25, 26], embedded systems [21], dynamical optimization [26], and systems biology [5, 17]. A hybrid ODE system consists of a family of ODEs, each of which is defined on a constraint set (i.e., invariant set) and forms a mode, with mode transitions occurring along state trajectories on boundaries of invariant sets following switching rules. As a distinct and intrinsic feature of hybrid dynamical systems, mode switching complicates various fundamental analytical, numerical, and design issues, for example, the issue of solution existence and uniqueness. Another critical issue, associated with short-time and finite-time hybrid dynamics, is whether infinitely many mode transitions exist in finite time, which is referred to as the Zeno behavior in the literature. More challenging, albeit practically important, issues pertain to long-time dynamic properties, e.g., stability, reachability/safety, and long-time system properties such as observability, which have found a wide range of applications [18].

Inspired by the study of global and long-time dynamics of hybrid systems, the present paper investigates positive invariance of an ODE system on a constraint set. Roughly speaking, a set is positively invariant if any system trajectory starting from the given set will remain in that set for all positive times; see Section 2 for a formal definition. The concept of positive invariance is essential in asymptotic analysis of unconstrained smooth dynamical systems and Lyapunov stability theory. For example, if a trajectory is contained in a compact set for all positive times, then its positive limit set is nonempty, compact, connected, and positively invariant [38, Proposition 1.1.14]. A

*Department of Mathematics and Statistics, University of Maryland Baltimore County, Baltimore, Maryland 21250, U.S.A. Email: shenj@umbc.edu.

prominent application of this result is in the proof of LaSalle’s invariance principle [20]. In the realm of hybrid systems, it has been recently recognized that positive invariance of each mode of a hybrid system plays a crucial role in addressing a variety of important long-time dynamic issues. For instance, it is shown in [34] that a positively invariant set of each linear mode of a class of linear hybrid systems can be used to characterize global switching behaviors of such a system. Positive invariance results also shed light on reachability analysis and safety verification with interesting applications in engineering [6] and systems biology [1, 5].

Given an ODE system Σ and a nonempty constraint set \mathcal{S} , two related but different positive invariance problems naturally arise:

- I Find conditions on Σ and \mathcal{S} such that \mathcal{S} is a positively invariant set of Σ ;
- II Characterize the (largest) positively invariant set of Σ contained in \mathcal{S} , where specific characterization issues include the existence of a positively invariant set and algebraic/geometric properties of such a set.

The first problem bears a larger body of the literature than the second, partly because it is more tractable and more relevant to control synthesis. Loosely speaking, this problem can be solved by imposing inward-pointing conditions on the vector field of Σ on the boundary of \mathcal{S} . Such conditions usually lead to verifiable algebraic results together with algebraic structure of \mathcal{S} . Some recent results, among many others, include [1] and [12] for positive invariance of a linear dynamics on a polyhedron and on a box respectively as well as [6] for a multi-affine dynamics on a rectangle. Also refer to [7] for the treatment of the similar problems and [37] for an application in control of multiple agents. The first positive invariance problem has been extended to differential inclusions with set-valued right-hand sides; see [4] and the related “viability theory”.

Being more analytic in nature, the second positive invariance problem, however, has received relatively less attention so far, in spite of equally important applications in hybrid systems, e.g., [18, 34]. One exception is [18, Theorem 3.1], which provides a neat necessary and sufficient condition for an affine dynamics on a polytope; the proof of this condition is based on an elegant use of the topological fixed point theorem and the convex structure of the problem. However, if the boundedness or stability-like assumption is removed, then characterization of the positively invariant set becomes rather nontrivial, both analytically and numerically, even for relatively simpler dynamics and constraint sets. Nonetheless numerous hybrid systems and many reachability analysis or safety verification problems possess unbounded constraint sets, which call for novel techniques to handle them.

In this paper, we address the positive invariance problem of the second kind with a focus on an affine dynamics on a (possibly unbounded) convex polyhedron. While such a problem can be formulated as an infinite dimensional linear programming and tackled from topological perspective, e.g., [16, Theorems 12-13] and [2], no finitely verifiable conditions are established. Combining the recent results for long-time linear dynamic analysis and the algebraic structure of lexicographical relation, we obtain verifiable necessary and sufficient existence conditions for positive invariance of an affine dynamics on a general polyhedron. It should be pointed out that although we concentrate on a seemingly simple dynamics, the analysis performed is far from trivial. Instead, it involves many nontrivial techniques from asymptotic analysis and lexicographic algebra, which have not been fully exploited to our best knowledge. By generalizing our recent treatment for a class of linear hybrid systems in [34], the positive invariance results are applied to obtain an explicit characterization of global switching properties of piecewise affine hybrid systems with each mode defined on a possibly unbounded polyhedral set. These characterization results are expected to lead to refined and less conservative stability conditions. Moreover, we exploit the positive invariance techniques developed in the paper to show decidability of exact safety verification problems for affine dynamics on semi-algebraic sets.

The rest of the paper is organized as follows. In Section 2, we formally define the positive invariance and safety verification problems and present preliminary technical results. In Section 3, a detailed development of the necessary and sufficient algebraic conditions for the existence of a positively invariant set on a polyhedron is given; special cases and extensions are discussed. Section 4 focuses on global switching properties of piecewise affine hybrid systems, for which verifiable algebraic conditions are established using the positive invariance results. In section 5, decidability analysis is performed for safety verification of affine systems with the aid of the long-time analysis techniques developed in Sections 2 and 3. The paper concludes with final remarks and discussions on future research directions in Section 6.

2 Problem Formulation and Preliminaries

2.1 Positive Invariance and Safety Verification

Consider an ODE system on \mathbb{R}^n

$$\dot{x} = f(x) \tag{1}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is assumed to be globally Lipschitz continuous. Let $x(t, x^0)$ denote the unique C^1 trajectory of (1) corresponding to an initial state $x^0 \in \mathbb{R}^n$.

Definition 1. [38] A set $\mathcal{N} \subseteq \mathbb{R}^n$ is called *positively invariant* (with respect to (1)) if the following implication holds: $x^0 \in \mathcal{N} \implies x(t, x^0) \in \mathcal{N}, \forall t \geq 0$. The set \mathcal{N} is called *invariant* if $x^0 \in \mathcal{N} \implies x(t, x^0) \in \mathcal{N}, \forall t \in \mathbb{R}$.

Two problems emerge from many applied dynamical systems modeled by (1):

- Positive invariance analysis (of the second kind): given a set \mathcal{S} in \mathbb{R}^n , characterize the subset $\mathcal{A} \subseteq \mathcal{S}$ that consists of all the x^0 satisfying the following implication

$$x^0 \in \mathcal{A} \implies x(t, x^0) \in \mathcal{S}, \forall t \geq 0. \tag{2}$$

- Safety verification on a time interval $\Delta \subseteq \mathbb{R}$: given two sets \mathcal{S}_0 and \mathcal{S}_f in \mathbb{R}^n such that $\mathcal{S}_0 \subseteq \mathcal{S}_f$, determine whether the following implication holds

$$x^0 \in \mathcal{S}_0 \implies x(t, x^0) \in \mathcal{S}_f, \forall t \in \Delta. \tag{3}$$

Typical time intervals include \mathbb{R} and $\mathbb{R}_+ \equiv [0, \infty)$.

The first problem bears the name *positive invariance* in that the subset \mathcal{A} is the largest positively invariant in \mathcal{S} with respect to (1) (provided that it is nonempty), due to the semi-group property of the trajectories. The set \mathcal{S}_f in the second problem represents the safe region of the system and the implication (3) means that any trajectory starting from the initial set \mathcal{S}_0 will not leave \mathcal{S}_f at any time in Δ , and thus remains *safe* from \mathcal{S}_0 on Δ . It is worth pointing out that the safety verification problem can be treated as a relevant positive invariance problem. For example, the dynamics is safe on the entire time domain (i.e. $\Delta = \mathbb{R}$) if and only if the initial set \mathcal{S}_0 is contained in $\mathcal{A} \cap \mathcal{A}'$, where \mathcal{A} and \mathcal{A}' are the largest positively invariant sets of the safe region \mathcal{S}_f with respect to $\dot{x} = f(x)$ and $\dot{x} = -f(x)$, respectively.

Before closing this subsection, we mention some basic facts for positive invariance analysis defined above. It is easy to see that if \mathcal{S} is closed and its positively invariant set \mathcal{A} given in (2) is nonempty, then \mathcal{A} is closed. Moreover, if f is affine and \mathcal{S} is convex, then \mathcal{A} is convex. Hence, if f is affine and \mathcal{S} is a closed convex polyhedron, then \mathcal{A} is closed and convex, albeit not necessarily polyhedral, upon its existence. If \mathcal{S} is additionally bounded (i.e., \mathcal{S} is a polytope), then \mathcal{A} is compact and convex.

2.2 Preliminary Technical Results

Two key techniques have been exploited in this paper to carry out the positive invariance analysis: long-time dynamic analysis of oscillatory dynamics and lexicographic relation. We present preliminary results associated with these techniques.

2.2.1 Long-time dynamic analysis of oscillatory modes

The following results provide major tools to deal with oscillatory modes (corresponding to complex eigenvalues of the defining matrix) in positive invariance analysis and safety verification of affine dynamics. These results are related to the so-called “almost periodic functions” [30].

Lemma 2. [34, Lemma 13] Given (finitely many) continuous periodic functions $g_i : \mathbb{R} \rightarrow [a_i, b_i]$ with frequency $\omega_i > 0$, where $[a_i, b_i] \subseteq \mathbb{R}$ and $i = 1, \dots, m$. Assume that each g_i is onto $[a_i, b_i]$ and the frequency ratio ω_i/ω_j is irrational for $i \neq j$. Then $(g_1(t), \dots, g_m(t))$ is dense on $[a_1, b_1] \times \dots \times [a_m, b_m]$, i.e., for any given $\tilde{y} \in [a_1, b_1] \times \dots \times [a_m, b_m]$ and any scalar $\varepsilon > 0$, there is a $\tilde{t} \geq 0$ such that $\|\tilde{y} - (g_1(\tilde{t}), \dots, g_m(\tilde{t}))\| \leq \varepsilon$.

The next result states that a nontrivial linear combination of sinusoidal functions has persistent sign alternating and its positive/negative variations are not diminishing as time goes.

Lemma 3. [34, Corollary 15] Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be $f(t) \equiv \sum_{i=1}^m [\alpha_i \cos(\omega_i t) + \beta_i \sin(\omega_i t)]$, where $\omega_i > 0$, $\omega_i \neq \omega_j$ whenever $i \neq j$, and $|\alpha_i| + |\beta_i| \neq 0$ for all i . Then there exist two scalars $\gamma_1 > 0$ and $\gamma_2 < 0$ such that for any t_* , two time instants $t_1, t_2 \in [t_*, \infty)$ exist satisfying $f(t_1) \geq \gamma_1$ and $f(t_2) \leq \gamma_2$.

In view of this result, we immediately have:

Corollary 4. Let $v^i \in \mathbb{R}^m$ and the scalars $\omega_i > 0$, θ_i be given, where $i = 1, \dots, \ell$ and $\omega_i \neq \omega_j$ whenever $i \neq j$. Then $\sum_{i=1}^{\ell} v^i \sin(\omega_i t + \theta_i) \geq 0$ for all $t \geq 0$ sufficiently large if and only if $v^i = 0$ for all $i = 1, \dots, \ell$.

We discuss more properties for the function $f(t) \equiv \sum_{i=1}^m [\alpha_i \cos(\omega_i t) + \beta_i \sin(\omega_i t)]$. For notational simplicity, let $d_i(t) \equiv \alpha_i \cos(\omega_i t) + \beta_i \sin(\omega_i t)$. By considering the rationality of ratios of the frequencies, we obtain the collection of (disjoint and distinct) equivalent classes $E_{\omega_j} = \{d_i(t) \mid \omega_i/\omega_j \text{ is rational}\}$ as shown in [34, Lemma 14]. Note that each equivalent class E_{ω_j} attains a basis frequency $\tilde{\omega}_s > 0$, namely, $\omega_i/\tilde{\omega}_s$ is a positive integer for any frequency ω_i associated with the function $d_i(t) \in E_{\omega_j}$. Let $E_{\tilde{\omega}_s}$ denote the equivalent class and let $q_{\tilde{\omega}_s}(t) \equiv \sum_{d_i \in E_{\tilde{\omega}_s}} d_i(t)$.

Then the following hold:

- (1) $q_{\tilde{\omega}_s}(\cdot)$ is a real-valued smooth and periodic function with the frequency $\tilde{\omega}_s$;
- (2) if $q_{\tilde{\omega}_s}(\cdot)$ is not identically zero, then it attains the maximal and minimal values $\sigma_{\tilde{\omega}_s} > 0$ and $\nu_{\tilde{\omega}_s} < 0$ on $(-\infty, \infty)$ respectively;
- (3) $q_{\tilde{\omega}_s}(\cdot)$ is onto $[\nu_{\tilde{\omega}_s}, \sigma_{\tilde{\omega}_s}]$;
- (4) the ratio of two basis frequencies associated with two distinct equivalent classes is irrational.

Suppose there are k equivalent classes $E_{\tilde{\omega}_s}$ and thus $f(t) = \sum_{s=1}^k q_{\tilde{\omega}_s}(t)$. Notice that while each $q_{\tilde{\omega}_s}$ is periodic, f is generally not and hence may not attain its maximum and minimum on \mathbb{R} .

In spite of this, the following lemma shows that its supremum (resp. infimum) is the sum of the maxima (resp. minima) of $q_{\tilde{\omega}_s}$'s. This observation is essential to decidability analysis of safety verification problems treated in Section 5.

Lemma 5. Let $\sigma_{\tilde{\omega}_s}$ and $\nu_{\tilde{\omega}_s}$ be defined above for the function f . Then $\sup_{[t_*, \infty)} f(t) = \sum_{s=1}^k \sigma_{\tilde{\omega}_s}$ and $\inf_{[t_*, \infty)} f(t) = \sum_{s=1}^k \nu_{\tilde{\omega}_s}$ for any $t_* \in \mathbb{R}$.

Proof. It is clear that $\sum_{s=1}^k \sigma_{\tilde{\omega}_s}$ and $\sum_{s=1}^k \nu_{\tilde{\omega}_s}$ are upper and lower bounds of f , respectively. To show $\sum_{s=1}^k \sigma_{\tilde{\omega}_s}$ is the least upper bound, we assume that each $q_{\tilde{\omega}_s}$ is not identically zero without loss of generality. Hence, by Lemma 2 and the properties (1-4) stated above, we see that for any $\varepsilon > 0$ sufficiently small, there is a $\tilde{t} \geq t_*$ such that $q_{\tilde{\omega}_s}(\tilde{t}) \in (\sigma_{\tilde{\omega}_s} - \varepsilon/k, \sigma_{\tilde{\omega}_s}]$ for all $s = 1, \dots, k$. Therefore, $(\sum_{s=1}^k \sigma_{\tilde{\omega}_s} - \varepsilon)$ is not an upper bound of f . Hence, $\sum_{s=1}^k \sigma_{\tilde{\omega}_s}$ is the least upper bound. Similarly, $\sum_{s=1}^k \nu_{\tilde{\omega}_s}$ is the greatest lower bound of f for any t_* . \square

It is worth noticing that Lemma 5 further implies: (1) $\sum_{s=1}^k \nu_{\tilde{\omega}_s} = \inf_{(-\infty, \infty)} f(t)$, and (2) $f(t) \geq \rho, \forall t$ for a scalar ρ if and only if $\sum_{s=1}^k \nu_{\tilde{\omega}_s} \geq \rho$.

2.2.2 Lexicographic relation

In this subsection, we discuss the lexicographical relation that is crucial to establish sufficient positive invariance conditions. An ordered real ℓ -tuple $a = (a_1, \dots, a_\ell)$ is called lexicographically nonnegative if either $a = 0$ or its first nonzero element (from the left) is positive and we write $a \succcurlyeq 0$. In the latter case, we call the first positive element the positive leading term/entry of the tuple. If a is not only lexicographically nonnegative but also nonzero, then a is called lexicographically positive and we write $a \succ 0$. For two tuples a and b , we write $a \succcurlyeq b$ if $(a - b) \succcurlyeq 0$. Hence, the lexicographical relation defines a linear order on the vector space of ℓ -tuples. An n -dimensional vector tuple (x^1, \dots, x^ℓ) is called lexicographically nonnegative (resp. positive) if each real tuple (x_i^1, \dots, x_i^ℓ) is lexicographically nonnegative (resp. positive) for all $i = 1, \dots, n$ and we write $(x^1, \dots, x^\ell) \succcurlyeq (\succ) 0$. The set of lexicographically nonnegative (resp. positive) real ℓ -tuples forms a convex, although not closed, cone in \mathbb{R}^ℓ . In what follows, we let \mathbb{R}_+^n (resp. \mathbb{R}_{++}^n) be the nonnegative (resp. positive) orthant of \mathbb{R}^n .

Lemma 6. Let $E : \mathbb{R}^n \rightarrow \mathbb{R}^{m_1}$ and $F : \mathbb{R}^n \rightarrow \mathbb{R}^{m_2}$ be two generalized positively homogeneous functions, i.e., $E_i(\lambda x) = p_i(\lambda)E_i(x)$, $F_j(\lambda x) = p_j(\lambda)F_j(x)$, $\forall x \in \mathbb{R}^n$, $\forall \alpha \in \mathbb{R}_+$, $\forall i, j$, where $p_i, p_j : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ are bijective and strictly increasing. Suppose there exist $x^1, \dots, x^k \in \mathbb{R}^n$ such that $(E(x^1), \dots, E(x^k)) \succ 0$ and $(F(x^1), \dots, F(x^k)) \succcurlyeq 0$. Then for any $a \in \mathbb{R}_{++}^{m_1}$, there exist nonnegative scalars μ_1, \dots, μ_k such that the positive leading terms $E_i(\mu_\ell x^\ell)$ and $F_s(\mu_p x^p)$ of any row of $(E(x^1), \dots, E(x^k))$ and $(F(x^1), \dots, F(x^k))$ satisfy $E_i(\mu_\ell x^\ell) \geq \left(\sum_{j=\ell+1}^k |E_i(\mu_j x^j)| \right) + a_i$ and

$$F_s(\mu_p x^p) \geq \left(\sum_{j=p+1}^k |F_s(\mu_j x^j)| \right), \text{ respectively.}$$

Proof. We assume, via suitable row switchings, that both $(E(x^1), \dots, E(x^k))$ and $(F(x^1), \dots, F(x^k))$ are of echelon form, namely, the nonzero row(s) are above all the zero rows and the positive leading term of a row is either at the same column or to the right of any leading term above it. We also assume that each column $z^i \equiv \begin{pmatrix} E(x^i) \\ F(x^i) \end{pmatrix} \in \mathbb{R}^{m_1+m_2}$ contains at least one positive leading term of some row, since otherwise we can simply drop such a column (or equivalently by choosing $\mu_i = 0$) without affecting the conclusion of the lemma. To find the desired μ_i 's, define $\theta_i \subseteq \{1, \dots, m_1+m_2\}$ to be the set of the indices of the positive leading terms in z^i for each $i = 1, \dots, k$. Clearly, θ_i is nonempty and $z_{\theta_i}^i > 0$ for each i . Let $\bar{a} \equiv \begin{pmatrix} a \\ 0 \end{pmatrix}$. We now start from the rightmost column, i.e., z^k , and choose $\mu_k := \max_{j \in \theta_k} p_j^{-1} \left(\frac{\bar{a}_j}{z_j^k} \right)$, where p_j^{-1} is the inverse function of p_j . It is clear that $p_j^{-1} : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is also bijective and strictly increasing. Hence, for any $j \in \theta_k$ corresponding to $E_j(x^k)$, we have

$$E_j(\mu_k x^k) = p_j(\mu_k) E_j(x^k) \geq \left[p_j \circ p_j^{-1} \left(\frac{\bar{a}_j}{z_j^k} \right) \right] E_j(x^k) \geq \bar{a}_j,$$

where we use $z_j^k \equiv E_j(x^k) > 0$ and the construction of μ_k . This also holds for $j \in \theta_k$ corresponding to $F_j(x^k)$. For $i = 1, \dots, k-1$, we then recursively compute μ_i (from $i = k-1$ to $i = 1$) as

$$\mu_i := \max_{j \in \theta_i} p_j^{-1} \left(\frac{\sum_{\ell=i+1}^k p_j(\mu_\ell) |z_j^\ell| + \bar{a}_j}{z_j^i} \right) \quad (4)$$

Notice that $\sum_{\ell=i+1}^k p_j(\mu_\ell) |z_j^\ell|$ equals to $\sum_{\ell=i+1}^k |E_j(\mu_\ell x^\ell)|$ or $\sum_{\ell=i+1}^k |F_j(\mu_\ell x^\ell)|$ for any $j \in \theta_i$. It is easy to verify via the echelon structure and an induction argument that these μ_i 's are the desired ones. \square

Corollary 7. Let E be the generalized positively homogeneous function defined above and suppose $(E(z^\ell), \dots, E(z^0)) \succ 0$ for $z^0, \dots, z^\ell \in \mathbb{R}^n$. Then there exist nonnegative scalars μ_1, \dots, μ_ℓ such that $E(z^0) + \sum_{i=1}^{\ell} E(\mu_i z^i) > 0$ and $E(z^j) + \sum_{i=j+1}^{\ell} E(\mu_{i-j} z^i) \geq 0$ for all $j = 1, \dots, \ell-1$. In particular, if $(z^\ell, \dots, z^0) \succ 0$ holds, then there exist nonnegative scalars μ_1, \dots, μ_ℓ such that $z^0 + \sum_{i=1}^{\ell} \mu_i z^i > 0$ and $z^j + \sum_{i=j+1}^{\ell} \mu_{i-j} z^i \geq 0$ for all $j = 1, \dots, \ell-1$.

Proof. Let $y^i \equiv E(z^i)$. Since $(y^\ell, \dots, y^0) \succ 0$, $(y^\ell, \dots, y^i) \succcurlyeq 0$ for all $i = 1, \dots, \ell$. Hence,

$$\begin{pmatrix} y^\ell & y^{\ell-1} & \dots & y^1 & y^0 \end{pmatrix} \succ 0$$

$$\begin{pmatrix} & y^\ell & \dots & y^2 & y^1 \\ & & \ddots & \vdots & \vdots \\ & & & y^\ell & y^{\ell-1} \\ & & & & y^\ell \end{pmatrix} \succcurlyeq 0$$

We also assume, via suitable row switchings, that the above tuple is of the echelon form as in the previous lemma. Letting $\mu_0 \equiv 1$, which corresponds to the rightmost column, and applying

Lemma 6, we obtain nonnegative reals μ_1, \dots, μ_ℓ such that the positive leading term y_j^s in the j th row of (y^ℓ, \dots, y^0) satisfies $E_j(\mu_s z^s) > \sum_{i=0}^{s-1} |E_j(\mu_i z^i)|$ and that the positive leading term y_k^s in the k th row of (y^ℓ, \dots, y^p) satisfies $E_k(\mu_{s-p} z^s) \geq \sum_{i=p}^{s-1} |E_k(\mu_{i-p} z^i)|$, where $p \in \{1, \dots, \ell - 1\}$.

Noting that $\mu_0 \equiv 1$ and $E_j(z^r) \equiv 0$ for all $r > s$, we have $E_j(z^0) + \sum_{i=1}^{\ell} E_j(\mu_i z^i) = E_j(\mu_s z^s) + \sum_{i=0}^{s-1} E_j(\mu_i z^i) > \sum_{i=0}^{s-1} |E_j(\mu_i z^i)| + \sum_{i=0}^{s-1} E_j(\mu_i z^i) \geq 0$. Hence, $E(z^0) + \sum_{i=1}^{\ell} E(\mu_i z^i) > 0$. Similarly, we can show $E(z^j) + \sum_{i=j+1}^{\ell} E(\mu_{i-j} z^i) \geq 0$ for all $j = 1, \dots, \ell - 1$. Finally, letting E be the identify mapping, we obtain the special case stated in the corollary. \square

3 Positive Invariance of Affine Dynamics on a Polyhedron

In this section, we are concerned with fundamental positive variance issues of an affine dynamics on a general convex polyhedron. In particular, we address the existence of a positively invariant set and finite verification of the existence conditions. Necessary and sufficient algebraic conditions are derived and extensions are discussed. Throughout this section, let the affine dynamics be $\dot{x} = Ax + d$ and a nonempty polyhedron be $\mathcal{P} = \{x \in \mathbb{R}^n \mid Cx \geq b\}$, where $A \in \mathbb{R}^{n \times n}$, $d \in \mathbb{R}^n$, $C \in \mathbb{R}^{m \times n}$, and $b \in \mathbb{R}^m$.

3.1 Positive Invariance of Linear Dynamics on a Polyhedron

To ease the presentation, we focus on linear dynamics first, i.e., $d = 0$. Without loss of generality, we assume throughout this section that the matrix A is of the real Jordan canonical form. Let J_{ij} be the j th real Jordan block associated with a possibly complex eigenvalue λ_i of A . Then A is the direct sum of J_{ij} 's and we write $A = \bigoplus_{i,j} J_{ij}$. For each J_{ij} , let C^{ij} denote the corresponding block in C . Let n_{ij} be the order of J_{ij} and $\bar{n}_i = \max_j n_{ij}$. For each real eigenvalue λ_i of A , suppose there are ℓ_i real Jordan blocks J_{ij} . For the j th Jordan block $J_{ij} \in \mathbb{R}^{n_{ij} \times n_{ij}}$, let the corresponding n -vector be of the form $v = (0, \dots, 0, (v^{ij})^T, 0, \dots, 0)^T$, where $v^{ij} \in \mathbb{R}^{n_{ij}}$. The collection of such the vectors, denoted by V_{ij} , is a subspace of \mathbb{R}^n isomorphic to $\mathbb{R}^{n_{ij}}$. The direct sum of all the subspaces V_{ij} 's is a subspace of \mathbb{R}^n given by $\{(0, \dots, 0, (v^{i1})^T, \dots, (v^{i\ell_i})^T, 0, \dots, 0)^T \in \mathbb{R}^n \mid v^{ij} \in \mathbb{R}^{n_{ij}}, j = 1, \dots, \ell_i\}$ and is isomorphic to $\mathbb{R}^{n_{i1} + \dots + n_{i\ell_i}}$, denoted by $V(\lambda_i) = \bigoplus_{j=1}^{\ell_i} \mathbb{R}^{n_{ij}}$. Each vector $v^i \in V(\lambda_i)$ is formed by vector stacking, i.e., $v^i = ((v^{i1})^T, \dots, (v^{i\ell_i})^T)^T$, where $v^{ij} \in \mathbb{R}^{n_{ij}}$. For notational convenience, we write $v^i = \bigoplus_j v^{ij}$. For a vector $u = (u_1, \dots, u_\ell)^T$ in \mathbb{R}^ℓ , we define the following lifting operator:

$$\mathcal{L}(u) := (u_2, \dots, u_\ell, 0)^T$$

Let \mathcal{L}^k be the composition of k -copies of \mathcal{L} , i.e. $\mathcal{L}^k = \underbrace{\mathcal{L} \circ \dots \circ \mathcal{L}}_{k\text{-times}}$. By convention, we let \mathcal{L}^0 be

the identity mapping, i.e., $\mathcal{L}^0(u) \equiv u$, and $\mathcal{L}^{-1}(u) \equiv 0, \forall u$. It is easy to verify that \mathcal{L} is a linear operator and satisfies $\mathcal{L}^{k_1} \circ \mathcal{L}^{k_2} = \mathcal{L}^{k_2} \circ \mathcal{L}^{k_1} = \mathcal{L}^{k_1+k_2}$ for any nonnegative integers k_1, k_2 . Moreover, for any $u \in \mathbb{R}^\ell$, $\mathcal{L}^k(u) \equiv 0$ if $k \geq \ell$. The lifting operator \mathcal{L} provides a compact way to express the solution of a linear dynamics defined by a Jordan canonical form. Indeed, for any Jordan block J_{ij} associated with a real eigenvalue λ_i and any vector $v^{ij} \in \mathbb{R}^{n_{ij}}$, we have

$$e^{J_{ij} t} v^{ij} = e^{\lambda_i t} \sum_{k=0}^{n_{ij}-1} \frac{t^k}{k!} \mathcal{L}^k(v^{ij})$$

Let $C^i = [C^{i1}, \dots, C^{i\ell_i}]$ and $J_i = \bigoplus_j J_{ij}$. Similarly, for any $v^i \equiv \bigoplus_j v^{ij} \in V(\lambda_i)$, we have

$$C^i e^{J_i t} v^i = \sum_j C^{ij} e^{J_{ij} t} v^{ij} = e^{\lambda_i t} \sum_{k=0}^{\bar{n}_i-1} \frac{t^k}{k!} \left(\sum_j C^{ij} \mathcal{L}^k(v^{ij}) \right), \quad (5)$$

where we recall $\bar{n}_i = \max_j n_{ij}$ and $\mathcal{L}^k(v^{ij}) = 0$ whenever $k \geq n_{ij}$. For such a $v^i \in V(\lambda_i)$ and an index set $\theta \subseteq \{1, \dots, m\}$ whose cardinality is denoted by $|\theta|$, we define the following $|\theta|$ -dimensional vector tuple

$$P(C_\theta^i, v^i) \equiv \left(\sum_j C_\theta^{ij} \mathcal{L}^{\bar{n}_i-1}(v^{ij}), \sum_j C_\theta^{ij} \mathcal{L}^{\bar{n}_i-2}(v^{ij}), \dots, \sum_j C_\theta^{ij} \mathcal{L}^0(v^{ij}) \right) \quad (6)$$

Proposition 8. For a pair (C^i, J_i) associated with a nonnegative real eigenvalue λ_i of A , if $v^i \equiv \bigoplus_j v^{ij} \in V(\lambda_i)$ and the index subsets $\theta, \theta' \subseteq \{1, \dots, m\}$ exist such that $P(C_\theta^i, v^i) \succ 0$ and $P(C_{\theta'}^i, v^i) \equiv 0$, then for any $a \in \mathbb{R}_{++}^{|\theta|}$, there exists $u^i \in V(\lambda_i)$ such that $C_\theta^i e^{J_i t} u^i \geq e^{\lambda_i t} a \geq a$ and $C_{\theta'}^i e^{J_i t} u^i \equiv 0, \forall t \geq 0$.

Proof. For notational convenience, let $z^k \equiv \sum_j C_\theta^{ij} \mathcal{L}^k(v^{ij})$ where $k = 0, \dots, \bar{n}_i - 1$. Therefore, we have $P(C_\theta^i, v^i) = (z^{\bar{n}_i-1}, \dots, z^0)$ such that $z^{\bar{n}_i-1} \geq 0$ and $(z^{\bar{n}_i-1}, \dots, z^0) \succ 0$. For the given $v^i = \bigoplus_j v^{ij}$, define $u^{ij} \equiv \sum_{s=0}^{\bar{n}_i-1} \mu_s \mathcal{L}^s(v^{ij})$, where $\mu_0 \equiv 1$ and the nonnegative real numbers $\mu_1, \dots, \mu_{\bar{n}_i-1}$ are to be determined. Let $u^i \equiv \bigoplus_j u^{ij} \in V(\lambda_i)$. Hence,

$$\begin{aligned} C_\theta^i e^{J_i t} u^i &= e^{\lambda_i t} \sum_{k=0}^{\bar{n}_i-1} \frac{t^k}{k!} \left(\sum_j C_\theta^{ij} \mathcal{L}^k(u^{ij}) \right) = e^{\lambda_i t} \sum_{k=0}^{\bar{n}_i-1} \frac{t^k}{k!} \left[\sum_j C_\theta^{ij} \left(\sum_{s=0}^{\bar{n}_i-1} \mu_s \mathcal{L}^{k+s}(v^{ij}) \right) \right] \\ &= e^{\lambda_i t} \sum_{k=0}^{\bar{n}_i-1} \frac{t^k}{k!} \left[\sum_j C_\theta^{ij} \left(\sum_{s=0}^{\bar{n}_i-1-k} \mu_s \mathcal{L}^{k+s}(v^{ij}) \right) \right] \\ &= e^{\lambda_i t} \sum_{k=0}^{\bar{n}_i-1} \frac{t^k}{k!} \left[\sum_{s=0}^{\bar{n}_i-1-k} \mu_s \sum_j C_\theta^{ij} \mathcal{L}^{k+s}(v^{ij}) \right] = e^{\lambda_i t} \sum_{k=0}^{\bar{n}_i-1} \frac{t^k}{k!} \left[\sum_{p=k}^{\bar{n}_i-1} \mu_{p-k} z^p \right] \end{aligned}$$

where $p \equiv k + s$, and we use the facts that the operator \mathcal{L} is linear and that $\mathcal{L}^{k+s}(u^{ij}) = 0$ if $k + s \geq \bar{n}_i - 1$ for all j . Furthermore, by using Corollary 7 and recalling $\mu_0 \equiv 1$ and $z^{\bar{n}_i-1} \geq 0$, we deduce that there exist nonnegative reals $\mu_p, p = 1, \dots, \bar{n}_i - 1$ such that $\sum_{p=0}^{\bar{n}_i-1} \mu_p z^p > 0$ and $\sum_{p=k}^{\bar{n}_i-1} \mu_{p-k} z^p \geq 0$

for all $k = 1, \dots, \bar{n}_i - 1$. In particular, we can further positively scale μ_p 's for each $p \geq 0$ such that $\sum_{p=0}^{\bar{n}_i-1} \mu_p z^p \geq a$. These results, together with $\lambda_i \geq 0$, yield $C_\theta^i e^{J_i t} u^i = e^{\lambda_i t} \sum_{k=0}^{\bar{n}_i-1} \frac{t^k}{k!} \left[\sum_{p=k}^{\bar{n}_i-1} \mu_{p-k} z^p \right] \geq$

$e^{\lambda_i t} \sum_{p=0}^{\bar{n}_i-1} \mu_p z^p \geq e^{\lambda_i t} a \geq a$ for all $t \geq 0$ as desired. Similarly, letting $\tilde{z}^p \equiv \sum_j C_{\theta'}^{ij} \mathcal{L}^p(v^{ij})$, we have

$$C_{\theta'}^i e^{J_i t} u^i = e^{\lambda_i t} \sum_{k=0}^{\bar{n}_i-1} \frac{t^k}{k!} \left[\sum_{p=k}^{\bar{n}_i-1} \mu_{p-k} \tilde{z}^p \right]. \text{ Since } P(C_{\theta'}^i, v^i) = (\tilde{z}^{\bar{n}_i-1}, \dots, \tilde{z}^0) \text{ and } P(C_{\theta'}^i, v^i) \equiv 0, \\ C_{\theta'}^i e^{J_i t} u^i \equiv 0 \text{ for all } t. \quad \square$$

Remark 9. For the vector tuple $(z^{\bar{n}_i-1}, \dots, z^0)$ in the above proof, let the index set ϕ consist of the indices corresponding to the positive leading terms of the tuple in the last column z^0 , i.e., $\phi \equiv \{j \mid z_j^0 \text{ is the positive leading term in the tuple}\}$. If ϕ is nonempty, then the lexicographical relation of z^p 's shows that $z_\phi^p = 0$ for all $p \geq 1$ and thus $C_\phi^i u^i = \sum_{p=0}^{\bar{n}_i-1} \mu_p z_\phi^p = \mu_0 z_\phi^0$. Moreover, if $z_\phi^0 \geq a_\phi > 0$, then it is easy to see that $\mu_0 \equiv 1$ needs not be positively scaled in order for $\sum_{p=0}^{\bar{n}_i-1} \mu_p z^p \geq a$ (but other μ_p 's may need). That is to say, μ_0 can always be chosen as one in this case. These observations will be used in the proof of Theorem 12. By convention, we assume that $P(C_\theta^i, v^i) \succ 0$ (resp. $\succ 0$) vacuously holds if the index set θ is empty.

The following example illustrates Proposition 8 and other related lexicographical conditions:

Example 10. Consider $J_i = \begin{bmatrix} \lambda_1 & 1 \\ 0 & \lambda_1 \end{bmatrix}$ with $\lambda_1 \geq 0$, $C_\theta^i = \begin{bmatrix} 0 & 5 \\ 2 & -1 \end{bmatrix}$, and $v^i = (-1, 1)^T$. Hence, $P(C_\theta^i, v^i) = (C_\theta^i \mathcal{L}(v^i), C_\theta^i v^i) = \begin{pmatrix} 0 & 5 \\ 2 & -3 \end{pmatrix} \succ 0$. Define $z^0 = (5, -3)^T$ and $z^1 = (0, 2)^T$. It follows from Corollary 7 that $\mu_1 = 2$ renders $z^0 + \mu_1 z^1 > 0$. Let $\mu_0 = 1$ and $u = \mu_1 \mathcal{L}(v^i) + \mu_0 v^i = (1, 1)^T$. Proposition 8 asserts that $C_\theta^i e^{J_i t} u = e^{\lambda_1 t} [\mu_0 z^1 t + (\mu_1 z^1 + \mu_0 z^0)] > 0, \forall t \geq 0$. Indeed, a straightforward computation shows that $C_\theta^i e^{J_i t} u = e^{\lambda_1 t} (5, 2t + 1)^T$. Moreover, since $e^{\lambda_1 t} [\mu_0 z^1 t + (\mu_1 z^1 + \mu_0 z^0)]$ is linear in (μ_0, μ_1) , one can positively scale (μ_0, μ_1) such that $e^{\lambda_1 t} [\mu_0 z^1 t + (\mu_1 z^1 + \mu_0 z^0)]$ is greater than any positive vector for all $t \geq 0$.

We introduce more notation. For the given vector $b = (b_1, \dots, b_m)^T \in \mathbb{R}^m$ characterizing the polyhedron \mathcal{P} , define three index sets: $\alpha \equiv \{i \mid b_i > 0\}$, $\beta \equiv \{i \mid b_i = 0\}$, and $\gamma \equiv \{i \mid b_i < 0\}$. For notational convenience, we let $b^+ = b_\alpha$, $b^0 = b_\beta = 0$, and $b^- = b_\gamma$. Likewise, we use C^+, C^0, C^- for the corresponding blocks $C_\alpha, C_\beta, C_\gamma$ in C respectively. Let $F, G : \mathbb{R} \rightarrow \mathbb{R}^\ell$ be two functions. We say that $F(t)$ tends to $G(t)$ as $t \rightarrow +\infty$ if for any $\varepsilon > 0$, there is $t_\varepsilon \geq 0$ such that $\|F(t) - G(t)\| \leq \varepsilon, \forall t \geq t_\varepsilon$. With this notion, we present the following proposition instrumental to the necessity proof of the main result.

Proposition 11. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ and $b_\ell \in \mathbb{R}$ be given such that $f(t) \geq b_\ell, \forall t \geq 0$. Suppose that $f(t)$ tends to $e^{\lambda t} \frac{t^p}{p!} \left[\rho_0 + \sum_{s=1}^k \rho_s \sin(\omega_s t + \theta_s) \right]$ as $t \rightarrow +\infty$, where the real tuple $(\rho_0, \rho_1, \dots, \rho_k) \neq 0$, p is a nonnegative integer, $\lambda \in \mathbb{R}_+$, $\omega_s \in \mathbb{R}_{++}$ with $\omega_i \neq \omega_j$ whenever $i \neq j$, and $\theta_s \in \mathbb{R}$. Then the following hold:

- (a) If $(\lambda, p) \succ (0, 0)$ or $b_\ell \geq 0$, then $\rho_0 > 0$;
- (b) If $\rho_0 \leq 0$, then $(\lambda, p) = (0, 0)$;
- (c) If $(\lambda, p) = (0, 0)$, then $\rho_0 \geq b_\ell$.

Proof. We prove (a) as follows. Let $h(t) \equiv \rho_0 + \sum_{s=1}^k \rho_s \sin(\omega_s t + \theta_s)$ and assume $\rho_0 \leq 0$. Since $(\rho_0, \rho_1, \dots, \rho_k) \neq 0$, either $\rho_0 < 0$ with $(\rho_1, \dots, \rho_k) = 0$ or $\rho_0 = 0$ with $(\rho_1, \dots, \rho_k) \neq 0$. For both the cases, we obtain a scalar $\eta < 0$, particularly via Lemma 3 for the latter case, such that for any $t_* \geq 0$, there exists $\tilde{t} \in [t_*, \infty)$ with $h(\tilde{t}) \leq \eta$. Consider the following two cases:

- (1) $(\lambda, p) \succ (0, 0)$. Since, as $t \rightarrow \infty$, $e^{\lambda t} t^p$ is arbitrarily large and $f(t)$ is sufficiently close to $e^{\lambda t} \frac{t^p}{p!} h(t)$, we deduce that for any $\zeta < 0$ and any small $\varepsilon > 0$, there exists $t' \geq 0$ such that $e^{\lambda t'} \frac{(t')^p}{p!} h(t') \leq \zeta$ and $|f(t') - e^{\lambda t'} \frac{(t')^p}{p!} h(t')| \leq \varepsilon$. This shows that for any $\zeta < 0$, there exists $t'' \geq 0$ with $f(t'') \leq \zeta/2$, which contradicts $f(t) \geq b_\ell, \forall t \geq 0$ for any given b_ℓ .

- (2) $b_\ell \geq 0$. We only need to look at the case $(\lambda, p) = (0, 0)$ since $(\lambda, p) \succ (0, 0)$ and the case $(\lambda, p) \succ (0, 0)$ has been treated in (1). Notice that $f(t)$ tends to $h(t)$ as $t \rightarrow +\infty$. Hence, via the property of $h(t)$ with $\eta < 0$ stated above, we see that for any $t_* \geq 0$, there exists $\tilde{t} \in [t_*, \infty)$ with $f(\tilde{t}) \leq \eta/2 < 0$, a contradiction.

Hence, we must have $\rho_0 > 0$ in (a). Statement (b) is a direct consequence of $(\lambda, p) \succ (0, 0)$ and (a). To show (c), notice that $(\lambda, p) = (0, 0)$ implies that $f(t)$ tends to $h(t)$ as $t \rightarrow +\infty$. If $(\rho_1, \dots, \rho_k) = 0$, then $\rho_0 \geq b_\ell$ holds true obviously. Otherwise, we apply Lemmas 3 and 5 to obtain $\rho_0 \geq b_\ell$. \square

The following result provides main necessary and sufficient existence conditions of a positively invariant set of a linear dynamics on a polyhedron; its extension to affine dynamics is given in Theorem 18. The proof for sufficiency relies on the lexicographic relation results developed before. The necessity proof shares the similar spirit in the recent paper [34]. However, a major difference is that [34] imposes observability-like conditions on the pairs (C^{ij}, J_{ij}) , while the present paper does not. The latter allows us to handle a general polyhedron on one hand, but considerably complicates the analysis on the other hand. This difficulty is overcome by introducing the lifting operator and making full advantage of the echelon structure of the lexicographic relation and the long-time dynamic analysis results as shown below. The proof is given in the Appendix in order to maintain a smooth paper flow.

Theorem 12. The positively invariant set \mathcal{A} of the linear dynamics on the polyhedron \mathcal{P} is nonempty if and only if the index set α is empty or there exist real eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_k \geq 0$ of A and $v^i \in V(\lambda_i)$, $i = 1, \dots, k$ such that

(a) $P^+ \equiv (P(C_\alpha^1, v^1), \dots, P(C_\alpha^k, v^k)) \succ 0$, $P^0 \equiv (P(C_\beta^1, v^1), \dots, P(C_\beta^k, v^k)) \succneq 0$, and

(b) the following implication holds:

$$P^- \equiv (P(C_\gamma^1, v^1), \dots, P(C_\gamma^k, v^k)) \not\leq 0 \implies \left[\lambda_k = 0, \sum_j C_\phi^{kj} v^{kj} \geq b_\phi, \sum_j C_\psi^{kj} v^{kj} \geq b_\psi, (P(C_{\gamma \setminus \psi}^1, v^1), \dots, P(C_{\gamma \setminus \psi}^k, v^k)) \succneq 0 \right], \quad (7)$$

where the index sets $\phi \subseteq \alpha$ and $\psi \subseteq \gamma$ consist of the indices corresponding to the rows whose positive (resp. negative) leading terms appear in the last columns of P^+ and P^- respectively.

We say a few more words about the implication (7). This condition should be read as: if P^- is not lexicographically nonnegative (namely, some row of P^- contains a negative leading term), then all the following four conditions must hold: (i) the last real eigenvalue λ_k is zero; (ii) each negative leading term of P^- only appears in the last column of P^- (which corresponds to the constant mode), and all the rest rows of P^- are lexicographically nonnegative; (iii) letting $\psi \subseteq \gamma$ denote the index set corresponding to the rows of P^- with negative leading terms (in the last column of P^-), then $\sum_j C_\psi^{kj} v^{kj} \geq b_\psi$; and (iv) letting $\phi \subseteq \alpha$ denote the index set corresponding to the rows of P^+ whose positive leading terms appear in the last column of P^+ , then $\sum_j C_\phi^{kj} v^{kj} \geq b_\phi$.

The conditions of Theorem 12 can be greatly simplified if A is a diagonal matrix. To see this, we assume, without loss of generality, that $A = \text{diag}(J_1, J_2, \dots, J_p)$, where the diagonal matrix block $J_i = \text{diag}(\lambda_i, \dots, \lambda_i)$, and the eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_p$. Accordingly, the $m \times n$ matrix C and an n -vector v can be partitioned as $C = [C^1 \ C^2 \ \dots \ C^p]$ and $v = ((v^1)^T, (v^2)^T, \dots, (v^p)^T)^T$ respectively, where C^i and v^i correspond to J_i . Note that for any J_i , each of its Jordan block is just the scalar λ_i so that $\bar{n}_i = 1$. Therefore, it follows from (6) that for any index set $\theta \subseteq \{1, \dots, m\}$, $P(C_\theta^i, v^i) = C_\theta^i v^i$. Hence P^+ in Theorem 12 becomes $(C_\alpha^1 v^1, \dots, C_\alpha^k v^k)$; the similar simplification can be made for P^0 and P^- . To further illustrate this discussion as well as the conditions of Theorem 12, we present the following example.

Example 13. Let $A \in \mathbb{R}^{4 \times 4}$, $C \in \mathbb{R}^{5 \times 4}$, and $b \in \mathbb{R}^5$ be

$$A = \begin{bmatrix} 5 & & & \\ & 4 & & \\ & & 0 & \\ & & & -6 \end{bmatrix}, \quad C = \begin{bmatrix} -1 & \star & \kappa & \star \\ 0 & 1 & 2 & \star \\ 0 & -1 & 1 & \star \\ 0 & -2 & 1 & \star \\ 0 & 0 & -0.5 & \star \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ 3 \\ 0.2 \\ 0 \\ -1 \end{bmatrix}$$

where \star denotes the uninteresting terms and κ is a real parameter we shall discuss below. Note that the polyhedron defined by C and b is unbounded, and the index sets $\alpha = \{1, 2, 3\}$, $\beta = \{4\}$, and $\gamma = \{5\}$. Let $v = (v^1, v^2, v^3, v^4)^T$, where each scalar v^i corresponds to the i th eigenvalue of A . Furthermore, C^i is the i th column of C in this case. Choose $v^1 = -1$ and $v^3 = 2$. We

obtain $P^+ = (C_\alpha^1 v^1, C_\alpha^3 v^3) = \begin{pmatrix} 1 & 2\kappa \\ 0 & 4 \\ 0 & 2 \end{pmatrix} \succ 0$, $P^0 = (C_\beta^1 v^1, C_\beta^3 v^3) = (0, 2) \succneq 0$, and $P^- = (C_\gamma^1 v^1, C_\gamma^3 v^3) = (0, -1) \not\prec 0$. To verify the implication (7), notice that v^3 corresponds to the zero eigenvalue, the negative leading term of P^- is in the last column of P^- , and the index sets $\phi = \{2, 3\}$ and $\psi = \{5\}$. Moreover, $C_\phi^3 v^3 = (4, 2)^T \geq b_\phi = (3, 0.2)^T$ and $C_\psi^3 v^3 = -1 \geq b_\psi = -1$. Hence, the positively invariant set \mathcal{A} exists. In fact, for any κ , $x = (-1 - 2|\kappa|, 0, 2, 0)^T \in \mathcal{A}$. Finally, we point out that a nonzero v^2 will not make $(C_\alpha^1 v^1, C_\alpha^2 v^2, C_\alpha^3 v^3) \succ 0$ since the leading terms of the second and third rows of C have opposite signs.

It is interesting to observe that conditions (a) and (b) of Theorem 12 do not involve the complex eigenvalues of A and thus can be cast as a constrained eigenvector problem. We elaborate more on this observation via the following example; computational issues pertaining to verification of the positive invariance conditions will be discussed in Remark 21.

Example 14. Let $A = \text{diag}(J_1, J_2) \in \mathbb{R}^{4 \times 4}$, $C \in \mathbb{R}^{3 \times 4}$, and $b \in \mathbb{R}^3$ be

$$J_1 = \begin{bmatrix} \lambda_1 & 1 \\ 0 & \lambda_1 \end{bmatrix}, \quad J_2 = \begin{bmatrix} \lambda_2 & \omega \\ -\omega & \lambda_2 \end{bmatrix}, \quad C = \begin{bmatrix} 0 & 5 & \star & \star \\ 2 & -1 & \star & \star \\ 0 & -1 & \star & \star \end{bmatrix}, \quad b = \begin{bmatrix} 3 \\ 0 \\ -4 \end{bmatrix}, \quad \text{with } \lambda_1 \geq 0, \omega > 0$$

Corresponding to J_1 and J_2 , we partition C and $v \in \mathbb{R}^4$ as $C = [C^1 \ C^2]$, and $v = ((v^1)^T, (v^2)^T)^T$, where $C^1 = (C_{\bullet 1}^1, C_{\bullet 2}^1)$ contains the first two columns of C and $v^1 = (v_1^1, v_2^1)^T \in \mathbb{R}^2$. Since J_1 has order 2, we obtain from (6) that, for any $\theta \subseteq \{1, 2, 3\}$, $P(C_\theta^1, v^1) = (C_{\theta 1}^1 v_2^1, C_{\theta 1}^1 v_1^1 + C_{\theta 2}^1 v_2^1)$. Moreover, the index sets $\alpha = \{1\}$, $\beta = \{2\}$, and $\gamma = \{3\}$. Since $P^+ = (0, 5 v_2^1)$, $P^0 = (2v_2^1, 2v_1^1 - v_2^1)$, and $P^- = (0, -v_2^1)$, we must have $v_2^1 > 0$ for $P^+ \succ 0$, but this implies $P^- \not\prec 0$. Therefore, if $\lambda_1 > 0$, then the condition (7) fails, which rules out the existence of the positively invariant set. On the other hand, if $\lambda_1 = 0$, then, as partially shown in Example 10, $v_1^1 = -1$ and $v_2^1 = 1$ satisfy the conditions of Theorem 12 and thus the positively invariant set exists.

Finally we informally explain why the oscillatory mode associated with the complex eigenvalues plays no role in the positive invariance conditions as follows. Consider a positively invariant trajectory $x(t, x^0)$. If $\lambda_2 > \lambda_1 \geq 0$, then the oscillatory mode would dominate as t sufficiently large. However, due to $\lambda_2 > 0$ and the persistent sign alternating property of the oscillatory mode (cf. Lemma 3), $Cx(t, x^0) \geq b$ cannot hold for large t , which contradicts the positive invariance of $x(t, x^0)$. Hence, such the oscillatory mode must vanish in $x(t, x^0)$. If $\lambda_2 < \lambda_1$, then as t is sufficiently large, the mode corresponding to J_1 will dominate. The case where $\lambda_2 = \lambda_1$ is subtle, but it can also be shown via the sign alternating property that the mode corresponding to J_1 plays a major role in $x(t, x^0)$. In all the cases, we obtain the (necessary) positive invariance conditions that only depend on the mode associated with λ_1 for sufficiently large t . On the other hand, under the positive invariance conditions, one can always use the lexicographical relation and perform positive scaling to obtain a positively invariant trajectory without considering the oscillatory mode.

3.2 Non-trivial Positive Invariance of Linear Dynamics on a Polyhedron

We discuss a special case in the positive invariance analysis of the linear dynamics. Recall that if the index set α is empty, then $0 \in \mathcal{A}$. We shall derive necessary and sufficient conditions for this case such that \mathcal{A} is nontrivial, i.e., $\mathcal{A} \neq \{0\}$, in this section. This result is useful to characterize global switching behaviors of piecewise affine systems to be discussed in Section 4. Notice that when \mathcal{P} is a polyhedral cone, i.e., $b = 0$, a similar result is obtained in [24] using the fixed point theorem. The interested reader may also refer to [7, 36] for a connection between this result and the Krein-Rutman Theorem in linear operator theory. The following results deal with the case where \mathcal{P} is other than a cone, i.e., γ is nonempty, using the dynamical systems approach developed in the previous section.

Lemma 15. Let (C^i, J_i) be a pair corresponding to a real eigenvalue λ_i of A . If there exists a nonzero $v^i \in V(\lambda_i)$ such that $P(C_\theta^i, v^i) \succcurlyeq 0$ for a nonempty index set θ , then there exists an eigenvector $u \in V(\lambda_i)$ of J_i (associated with λ_i) such that $C_\theta^i u \geq 0$.

Proof. Given the nonzero $v^i \equiv \bigoplus_j v^{ij} \in V(\lambda_i)$, let ℓ be the largest integer such that $\mathcal{L}^\ell(v^{ij}) \neq 0$ for some j , namely, $\mathcal{L}^p(v^{ij}) = 0, \forall j$ if $p > \ell$. Since $v^i \neq 0$, it is clear that ℓ exists and $0 \leq \ell \leq \bar{n}_i - 1$. For this ℓ , define the index set $\vartheta \equiv \{j \mid \mathcal{L}^\ell(v^{ij}) \neq 0\}$ which is nonempty. By the property of the lifting operator \mathcal{L} discussed in the previous section, we have, for any $j \in \vartheta$, $\mathcal{L}^\ell(v^{ij}) = (\rho_j, 0, \dots, 0)^T$, where $\rho_j \neq 0$. Therefore $\sum_j C_\theta^{ij} \mathcal{L}^\ell(v^{ij}) = \sum_{j \in \vartheta} \rho_j C_{\theta 1}^{ij}$, where $C_{\theta 1}^{ij}$ denotes the first column of C_θ^{ij} .

This thus shows

$$P(C_\theta^i, v^i) = (0, \dots, 0, \sum_j C_\theta^{ij} \mathcal{L}^\ell(v^{ij}), \star, \dots, \star) = (0, \dots, 0, \sum_{j \in \vartheta} \rho_j C_{\theta 1}^{ij}, \star, \dots, \star) \succcurlyeq 0$$

where \star denotes the terms we are not interested in. Consequently, $\sum_{j \in \vartheta} \rho_j C_{\theta 1}^{ij} \geq 0$. It is also noted that for each $j \in \vartheta$, $\mathcal{L}^\ell(v^{ij})$ is an eigenvector of J_{ij} associated with λ_i . Hence, $u = \bigoplus_j \mathcal{L}^\ell(v^{ij})$ is an eigenvector of J_i associated with λ_i . The lemma thus holds as $C_\theta^i u = \sum_{j \in \vartheta} \rho_j C_{\theta 1}^{ij} \geq 0$. \square

Theorem 16. Consider the linear dynamics and the polyhedron \mathcal{P} with $\alpha = \emptyset$, $C^0 \equiv C_\beta$ and $C^- \equiv C_\gamma$. Then \mathcal{A} is nontrivial if and only if any one of the following conditions holds:

- (a) (C, A) is an unobservable pair, i.e., there exists a nonzero vector $v \in \overline{O}(C, A)$;
- (b) there exists an eigenvector v associated with a real eigenvalue λ of A such that $C^0 v \geq 0$ and the following implication holds: $\lambda > 0 \implies C^- v \geq 0$;
- (c) there exist a complex eigenvalue $(\lambda + \iota\omega)$ of A (with $\omega \neq 0$) and an associated complex eigenvector $(u + \iota v)$ such that $\lambda \leq 0$ and $C^0 u = C^0 v = 0$.

Proof. ‘‘Sufficiency’’. Case (a) is trivial. We consider (b) and (c) as follow:

Case (b). Let λ be a real eigenvalue of A and $v \neq 0$ be an associated eigenvector satisfying the conditions stated in (b). First, it is easy to see $C^0 e^{At} v = e^{\lambda t} C^0 v \geq 0, \forall t \geq 0$. Consider two subcases: (i) $\lambda > 0$; and (ii) $\lambda \leq 0$. For the first subcase, we have $C^- e^{At} v = e^{\lambda t} C^- v \geq 0 \geq b^-, \forall t \geq 0$. Thus $v \in \mathcal{A}$. For subcase (ii) where $\lambda \leq 0$, we have a scalar $\varepsilon > 0$ such that $C^-(\varepsilon v) \geq b^-$. Notice that $C^0 e^{At}(\varepsilon v) \geq 0, \forall t \geq 0$. Define the index set $\phi = \{i \mid (C^- v)_i \geq 0\}$. Clearly, $[C^-(\varepsilon v)]_\phi \geq 0$ and $0 > [C^-(\varepsilon v)]_{\bar{\phi}} \geq (b^-)_{\bar{\phi}}$. Therefore, $[C^- e^{At}(\varepsilon v)]_\phi = e^{\lambda t} [C^-(\varepsilon v)]_\phi \geq 0 \geq (b^-)_\phi$ and $[C^- e^{At}(\varepsilon v)]_{\bar{\phi}} = e^{\lambda t} [C^-(\varepsilon v)]_{\bar{\phi}} \geq [C^-(\varepsilon v)]_{\bar{\phi}} \geq (b^-)_{\bar{\phi}}$ for all $t \geq 0$ (because $\lambda \leq 0$).

In conclusion, $C^-e^{At}(\varepsilon v) \geq b^-, \forall t \geq 0$. Thus $0 \neq \varepsilon v \in \mathcal{A}$.

Case (c). Notice that $v \neq 0$ and the condition $C^0u = C^0v = 0$ implies $(u + \iota v) \in \overline{O}(C^0, A)$, i.e., $C^0e^{At}(u + \iota v) = 0, \forall t$. Therefore, $C^0e^{At}v = 0, \forall t \geq 0$. Moreover, observing that $C^-e^{At}v = e^{\lambda t} [\cos(\omega t)C^-v + \sin(\omega t)C^-u]$ and $\lambda \leq 0$, we deduce that $\|C^-e^{At}v\|_2 \leq \delta, \forall t \geq 0$ for some scalar $\delta > 0$. Hence, there exists a scalar $\varepsilon > 0$ such that $\|C^-e^{At}(\varepsilon v)\|_2 \leq \min_{i \in \gamma}(|b_i^-|), \forall t \geq 0$. Therefore $C^-e^{At}(\varepsilon v) \geq b^-$ and $C^0e^{At}(\varepsilon v) \equiv 0$ for all $t \geq 0$. This implies $\varepsilon v \in \mathcal{A}$.

“Necessity”. Let $0 \neq x^* \in \mathcal{A}$, i.e., $Ce^{At}x^* \geq b, \forall t \geq 0$. To avoid triviality, we assume $x^* \notin \overline{O}(C, A)$. Since $Ce^{At}x^*$ is not identically zero, there is an eigenvalue λ_i associated with the largest non-vanishing mode in $Ce^{At}x^*$ as shown in the proof of Theorem 12. We consider two cases as follows:

(N1) $\lambda_i > 0$. In this case, by the similar argument in the necessity proof of Theorem 12, we deduce that there exists a nonzero $v^i \in V(\lambda_i)$ such that $P(C^i, v^i) \succ 0$. This thus yields, via Lemma 15, an eigenvector u of J_i such that $C^i u \geq 0$. Appropriately expanding u to a vector $v \in \mathbb{R}^n$ by adding the zero subvectors, we obtain an eigenvector v associated with λ_i such that $C^0v \geq 0$ and $C^-v \geq 0$. This leads to (b).

(N2) $\lambda_i \leq 0$. We consider two subcases: (i) $x^* \notin \overline{O}(C^0, A)$; or (ii) $x^* \in \overline{O}(C^0, A)$ but $x^* \notin \overline{O}(C^-, A)$. For the first subcase, $C^0e^{At}x^* = C_\beta e^{At}x^* \geq 0, \forall t \geq 0$ but is not identically zero. Hence, $P(C_\beta^i, v^i) \succ 0$ for some $0 \neq v^i \in V(\lambda_i)$. Following the argument in (N1), we thus obtain an eigenvector u associated with λ_i such that $C_\beta u = C^0 u \geq 0$. This results in (b) by observing that the implication in (b) vacuously holds as $\lambda_i \leq 0$. For the second subcase, we have the non-vanishing terms in $C_\gamma e^{At}x^*$ corresponding to λ_i as

$$e^{\lambda_i t} \sum_{k=0}^{\bar{n}_i-1} \frac{t^k}{k!} \left(\sum_j C_\gamma^{ij} \mathcal{L}^k(v^{ij}) + \sum_s (u_s^{ik})_\gamma \sin(\omega_s^{ik} t + \theta_s^{ik}) \right)$$

where $v^{ij}, u_s^{ik}, \omega_s^{ik}, \theta_s^{ik}$ depend on C, A, x^* only. If all $(u_s^{ik})_\gamma = 0$, then $v^i = \bigoplus_j v^{ij} \neq 0$. But since $C_\beta^i e^{J_i t} v^i \equiv 0, \forall t$, we have $P(C_\beta^i, v^i) \equiv 0$ which leads to an eigenvector v of A associated with λ_i such that $C_\beta v \equiv C^0 v \geq 0$. This gives rise to (b). If, however, $(u_s^{ik})_\gamma \neq 0$ for some s and k , then $u_s^{ik} \neq 0$ and we must have a complex eigenvector $(u + \iota v)$ associated with the complex eigenvalue $\lambda_i + \iota \omega_s^{ik}$ such that $C^0(u + \iota v) = 0$, or equivalently $C^0u = C^0v = 0$. Since $\lambda_i \leq 0$, we obtain (c) as desired. \square

It is worth pointing out that the algebraic conditions in the above theorem are finitely verifiable.

3.3 Positive Invariance of Affine Dynamics on a Polyhedron

We now return to positive invariance analysis of the affine dynamics $\dot{x} = Ax + d$ on the convex polyhedron $\mathcal{P} = \{x \in \mathbb{R}^n \mid Cx \geq b\}$, where $A \in \mathbb{R}^{n \times n}$ and $d \in \mathbb{R}^n$. For the given A and d , we can uniquely decompose d into $d = d_c + d_n$, where d_c is the orthogonal projection of d onto the column space of A and d_n is the projection onto the null space of A^T such that d_c is orthogonal to d_n . Since $d_c = Au_c$ for some vector u_c , the state transformation $\tilde{x} = x + u_c$ converts the affine dynamics into $\dot{\tilde{x}} = A\tilde{x} + d_n$ and the polyhedron into $\tilde{\mathcal{P}} = \{\tilde{x} \mid C\tilde{x} \geq \tilde{b}\}$, where $\tilde{b} \equiv b + Cu_c$. Hence, if $d_n = 0$, namely, d is in the range of A , then the affine dynamics can be transformed into the linear dynamics such that the results in the previous section follow. Consequently, we assume, without loss of generality, that $d \neq 0$ is in the null space of A^T . Equivalently, this assumes that d is an eigenvector of A^T associated with the zero eigenvalue which further implies that A has the zero eigenvalue.

To develop the necessary and sufficient condition for the existence of the positively invariant set of the affine dynamics, we present a technical lemma extended from Proposition 8. Let $J_i = \bigoplus_j J_{ij}$ be the Jordan block associated with the zero eigenvalue λ_i and $d^i = \bigoplus_j d^{ij}$ be an eigenvector of J_i^T .

Hence, $d^{ij} = (0, \dots, 0, d_{n_{ij}}^{ij})^T$, where $d_{n_{ij}}^{ij} \neq 0$ is the last element of d^{ij} . Recall that $\bar{n}_i = \max_j n_{ij}$ where n_{ij} is the order of J_{ij} . For a given index set θ and $v^i \in V(0)$, define the vector tuple:

$$P(C_\theta^i, v^i, d^i) \equiv \left(\sum_j C_\theta^{ij} [\mathcal{L}^{\bar{n}_i}(v^{ij}) + \mathcal{L}^{\bar{n}_i-1}(d^{ij})], \dots, \sum_j C_\theta^{ij} [\mathcal{L}(v^{ij}) + \mathcal{L}^0(d^{ij})], \sum_j C_\theta^{ij} \mathcal{L}^0(v^{ij}) \right) \quad (8)$$

Lemma 17. Let d^i be an eigenvector of J_i^T associated with the eigenvalue $\lambda_i = 0$. For a given vector $b \in \mathbb{R}^m$, if there exist $v^i \in V(\lambda_i)$ and the index subsets $\theta, \theta' \subseteq \{1, \dots, m\}$ such that $P(C_\theta^i, v^i, d^i) \succ 0$, $P(C_{\theta'}^i, v^i, d^i) \equiv 0$, and $\sum_j C_\phi^{ij} \mathcal{L}^0(v^{ij}) \geq b_\phi$, where $\phi \subseteq \theta$ consists of the indices corresponding to the positive leading term of the tuple in the last column of $P(C_\theta^i, v^i, d^i)$, then there exists $u^i \in V(\lambda_i)$ such that $C_\theta^i [e^{J_i t} u^i + \int_0^t e^{J_i(t-\tau)} d\tau d^i] \geq b_\theta$, $C_{\theta'}^i [e^{J_i t} u^i + \int_0^t e^{J_i(t-\tau)} d\tau d^i] \equiv 0$, $\forall t \geq 0$.

Proof. Let $d^i = \bigoplus_j d^{ij}$, where d^{ij} is in the null space of the Jordan block J_{ij}^T associated with the zero eigenvalue λ_i and at least one of d^{ij} 's is nonzero. Since $d^{ij} = (0, \dots, 0, d_{n_{ij}}^{ij})^T$ with $d_{n_{ij}}^{ij} \neq 0$, we have $\int_0^t e^{J_{ij}(t-\tau)} d\tau d^{ij} = \sum_{k=0}^{n_{ij}-1} \frac{t^{k+1}}{(k+1)!} (J_{ij})^k d^{ij} = \sum_{k=0}^{n_{ij}-1} \frac{t^{k+1}}{(k+1)!} \mathcal{L}^k(d^{ij})$. Let $z^0 \equiv \sum_j C_\theta^{ij} \mathcal{L}^0(v^{ij})$ and $z^k \equiv \sum_j C_\theta^{ij} [\mathcal{L}^k(v^{ij}) + \mathcal{L}^{k-1}(d^{ij})]$ for $k = 1, \dots, \bar{n}_i$. Likewise, let $y^0 \equiv \sum_j C_{\theta'}^{ij} \mathcal{L}^0(v^{ij})$ and $y^k \equiv \sum_j C_{\theta'}^{ij} [\mathcal{L}^k(v^{ij}) + \mathcal{L}^{k-1}(d^{ij})]$ for $k = 1, \dots, \bar{n}_i$. Hence, $(z^{\bar{n}_i}, \dots, z^0) \succ 0$ with $z_\phi^0 \geq b_\phi$ and $(y^{\bar{n}_i}, \dots, y^0) \equiv 0$. By Corollary 7 and Remark 9, we obtain nonnegative scalars $\mu_1, \dots, \mu_{\bar{n}_i}$ with $\mu_0 \equiv 1$ such that $\sum_{s=0}^{\bar{n}_i} \mu_s z^s \geq b_\theta$ and $\sum_{s=k}^{\bar{n}_i} \mu_{s-k} z^s \geq 0$ for all $k = 1, \dots, \bar{n}_i$, and $\sum_{s=k}^{\bar{n}_i} \mu_{s-k} y^s \equiv 0$ for all $k \geq 0$. Let $u^{ij} \equiv \mu_0 \mathcal{L}^0(v^{ij}) + \sum_{s=1}^{\bar{n}_i} \mu_s [\mathcal{L}^s(v^{ij}) + \mathcal{L}^{s-1}(d^{ij})]$, and $u^i \equiv \bigoplus_j u^{ij} \in V(\lambda_i)$. Hence,

$$\begin{aligned} & C_\theta^i [e^{J_i t} u^i + \int_0^t e^{J_i(t-\tau)} d\tau d^i] \\ &= \sum_{k=0}^{\bar{n}_i} \frac{t^k}{k!} \left(\sum_j C_\theta^{ij} [\mathcal{L}^k(u^{ij}) + \mathcal{L}^{k-1}(d^{ij})] \right) \\ &= \sum_{k=0}^{\bar{n}_i} \frac{t^k}{k!} \left(\sum_j C_\theta^{ij} \left[\sum_{s=1}^{\bar{n}_i-k} \mu_s [\mathcal{L}^{k+s}(v^{ij}) + \mathcal{L}^{k+s-1}(d^{ij})] + \mu_0 \mathcal{L}^k(v^{ij}) + \mu_0 \mathcal{L}^{k-1}(d^{ij}) \right] \right) \\ &= \sum_{k=0}^{\bar{n}_i} \frac{t^k}{k!} \left(\sum_j C_\theta^{ij} \sum_{s=0}^{\bar{n}_i-k} \mu_s [\mathcal{L}^{k+s}(v^{ij}) + \mathcal{L}^{k+s-1}(d^{ij})] \right) \\ &= \sum_{k=0}^{\bar{n}_i} \frac{t^k}{k!} \left(\sum_{p=k}^{\bar{n}_i} \mu_{p-k} \sum_j C_\theta^{ij} [\mathcal{L}^p(v^{ij}) + \mathcal{L}^{p-1}(d^{ij})] \right) = \sum_{k=0}^{\bar{n}_i} \frac{t^k}{k!} \left[\sum_{p=k}^{\bar{n}_i} \mu_{p-k} z^p \right], \end{aligned}$$

where we recall $\mathcal{L}^{-1}(w) \equiv 0$ for any w by convention. Therefore, we have $C_\theta^i [e^{J_i t} u^i + \int_0^t e^{J_i(t-\tau)} d\tau d^i] \geq \sum_{p=0}^{\bar{n}_i} \mu_p z^p \geq b_\theta$, $\forall t \geq 0$. Similarly, we obtain

$$C_{\theta'}^i [e^{J_i t} u^i + \int_0^t e^{J_i(t-\tau)} d\tau d^i] = \sum_{k=0}^{\bar{n}_i} \frac{t^k}{k!} \left[\sum_{p=k}^{\bar{n}_i} \mu_{p-k} y^p \right] = 0, \quad \forall t \geq 0. \quad \square$$

Theorem 18. Consider the affine dynamics $\dot{x} = Ax + d$, where $d \neq 0$ is in the null space of A^T . The positively invariant set \mathcal{A} of the affine dynamics on \mathcal{P} is nonempty if and only if there exist real eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_{k-1} > 0 = \lambda_k$ of A and $v^i \in V(\lambda_i)$, $i = 1, \dots, k$ such that $(P(C_\alpha^1, v^1), \dots, P(C_\alpha^{k-1}, v^{k-1}), P(C_\alpha^k, v^k, d^k)) \succ 0$ with $\sum_j C_\phi^{kj} v^{kj} \geq b_\phi$, $(P(C_\beta^1, v^1), \dots, P(C_\beta^{k-1}, v^{k-1}), P(C_\beta^k, v^k, d^k)) \succcurlyeq 0$, and the following implication holds:

$$(P(C_\gamma^1, v^1), \dots, P(C_\gamma^{k-1}, v^{k-1}), P(C_\gamma^k, v^k, d^k)) \not\prec 0 \implies \\ [\sum_j C_\psi^{kj} v^{kj} \geq b_\psi, (P(C_{\gamma \setminus \psi}^1, v^1), \dots, P(C_{\gamma \setminus \psi}^k, v^k, d^k)) \succcurlyeq 0],$$

where the index sets $\phi \subseteq \alpha$ and $\psi \subseteq \gamma$ consist of the indices corresponding to the rows whose positive (resp. negative) leading terms appear in the last columns of the respective tuples.

Proof. Since $d \neq 0$ is in the null space of A^T , A has the zero eigenvalue and thus A always has a nonnegative eigenvalue. The sufficiency of the theorem can be proved in the similar fashion via the corresponding argument in Theorem 12 by making use of Lemma 17 and $\lambda_k = 0$ for the zero eigenvalue mode. To prove the necessity, let $x^* \in \mathcal{A}$. It is clear that

$$C[e^{At}x^* + \int_0^t e^{A(t-\tau)}d\tau d] = \sum_{\lambda_i \neq 0} \left\{ e^{\lambda_i t} \sum_{k=0}^{\hat{n}_i-1} \frac{t^k}{k!} \left(\sum_j C^{ij} \mathcal{L}^k(v^{ij}) + \sum_s u_s^{ik} \sin(\omega_s^{ik}t + \theta_s^{ik}) \right) \right\} \\ + \sum_{k=0}^{\hat{n}_0} \frac{t^k}{k!} \left(\sum_j C^{0j} [\mathcal{L}^k(v^{0j}) + \mathcal{L}^{k-1}(d^{0j})] + \sum_s u_s^{0k} \sin(\omega_s^{0k}t + \theta_s^{0k}) \right),$$

where λ_i and \hat{n}_i are defined in the same way as in Theorem 12, v^{0j} and u_s^{0k} correspond to the zero eigenvalue and (possibly existing) imaginary eigenvalues of A respectively, which depend on C, A, x^* only. Let the real parts of the eigenvalues of A be labeled in a descending order, i.e., $\lambda_1 > \lambda_2 > \dots > \lambda_{k-1} > \lambda_k = 0 > \lambda_{k+1} > \dots > \lambda_q$. Define the vector tuples $P(C^i, v^i)$ corresponding to positive λ_i , $i = 1, \dots, k-1$ in the same fashion as in Theorem 12. In particular, if λ_i corresponds to a strictly complex eigenvalue of A , then λ_i must be positive and $P(C^i, v^i) \equiv 0$. We also define $P(C^k, v^k, d^k)$ for the zero eigenvalue λ_k with $v^k = \bigoplus_j v^{0j} \in V(\lambda_k)$ and $d^k = \bigoplus_j d^{0j}$, which is an eigenvector of J_0^T . The remaining proof thus follows from the necessity arguments in (N1–N3) of Theorem 12 by considering the large-time dominating mode of $C_\ell[e^{At}x^* + dt]$ for each $\ell \in \alpha$, $\ell \in \beta$ and $\ell \in \gamma$ respectively. Recall that if any index set is empty, then the associated lexicographical relation is assumed to hold vacuously; see Remark 9. Lastly, since A has the zero eigenvalue, by removing the positive λ_i that corresponds to a strictly complex eigenvalue of A and its corresponding zero block in the obtained tuples, we attain the desired nonnegative real eigenvalues $\tilde{\lambda}_i$ and associated subvectors $v^i \in V(\tilde{\lambda}_i)$ as explained at the end of Theorem 12. \square

We present the following example to illustrate the above conditions.

Example 19. Consider Example 13, where A is diagonal, and $d = (0, 0, \rho, 0)^T$ with $\rho \neq 0$ such that d is in the null space of A^T . Hence, $d^i = 0$, $i = 1, 2, 4$ and $d^3 = \rho$. Recall that the polyhedron is unbounded and $\alpha = \{1, 2, 3\}$, $\beta = \{4\}$, and $\gamma = \{5\}$. Since d^3 corresponds to the zero eigenvalue, it follows from (8) that $P(C_\theta^3, v^3, d^3) = (C_\theta^3 d^3, C_\theta^3 v^3)$ for any index set $\theta \subseteq \{1, \dots, 5\}$. Let $\rho > 0$. Note that no matter what v^1 , v^2 and v^3 are chosen, $(P(C_\gamma^1, v^1), P(C_\gamma^2, v^2), P(C_\gamma^3, v^3, d^3)) = (0, 0, -0.5\rho, -0.5v^3) \not\prec 0$ and the negative leading term -0.5ρ does not appear in the last column. Hence, the conditions of Theorem 18 fail such that the positively invariant set does not exist. This result can be directly verified as the last entry of $C[e^{At}x + \int_0^t e^{A(t-\tau)}d\tau d]$ takes the form $-0.5(\rho t + x_3 + x_4 e^{-6t})$, which cannot be greater than -1 for all $t \geq 0$. Now consider $\rho < 0$. Since

$$(P(C_\alpha^1, v^1), P(C_\alpha^2, v^2), P(C_\alpha^3, v^3, d^3)) = \begin{pmatrix} -v^1 & \star & \kappa\rho & \kappa v^3 \\ 0 & v^2 & 2\rho & 2v^3 \\ 0 & -v^2 & \rho & v^3 \end{pmatrix},$$

no v^1, v^2, v^3 will make the above vector tuple lexicographically positive. Hence, the positively invariant set does not exist either. Finally, we comment that if $\rho > 0$ and we change the (5, 3)-entry of C from -0.5 to a nonnegative number, then the conditions of Theorem 18 are satisfied for any κ , which ascertains the existence of the positively invariant set.

Remark 20. An interesting special case is when the convex polyhedron \mathcal{P} is bounded, i.e., \mathcal{P} is a polytope. For a general affine dynamics $\dot{x} = Ax + d$ with (possibly zero) d in the null space of A^T , it has been shown in [18, Theorem 3.1] via the fixed point argument that the associated positively invariant set \mathcal{A} exists if and only if $Av + d = 0$ for some $v \in \mathcal{P}$. The long-time dynamic analysis techniques in Theorems 12 and 18 provide an alternative proof for the necessity of this result (the sufficiency is trivial). Indeed, letting $x^* \in \mathcal{A}$ and noticing $e^{At}x^* + \int_0^t e^{A(t-\tau)}d\tau d$ is

bounded on $[0, \infty)$, we have $e^{At}x^* + \int_0^t e^{A(t-\tau)}d\tau d = \bigoplus_j [v^{0j} + t(\mathcal{L}(v^{0j}) + d^{0j}) + h_j(t)]$, where v^{0j} corresponds to the zero eigenvalue, $d = \bigoplus_j d^{0j}$, and $h_j(t)$ contains all the terms associated with the eigenvalues with non-positive real parts. Hence $\mathcal{L}(v^{0j}) + d^{0j} = 0$ for all j due to the boundedness. By (c) of Proposition 11, we deduce $\sum_j C^{0j} v^{0j} \geq b$. Moreover, $\mathcal{L}(v^{0j}) = J_{0j} v^{0j}$ holds. Hence, by appropriately expanding v^0 , we obtain v satisfying $Cv = \sum_j C^{0j} v^{0j} \geq b$ and $Av = \bigoplus_j J_{0j} v^{0j} = -\bigoplus_j d^{0j} = -d$. This shows that $v \in \mathcal{P}$ and $Av + d = 0$ as expected. This result can also be obtained by directly checking the conditions of Theorem 18 with the observations that $v^i = 0$ for $\lambda_i > 0$ and $\phi = \alpha$ if $\alpha \neq \emptyset$ due to the boundedness of \mathcal{P} .

Remark 21. The necessary and sufficient conditions presented in Theorems 12 and 18 pose the question of whether they can be verified by a finite procedure, or namely whether the verification problem is decidable. Upon knowing the eigenvalues of the matrix A , the answer is affirmative since the matrix A has finitely many nonnegative eigenvalues and checking the algebraic conditions in the theorems, especially the lexicographical relation, can be formulated as a feasibility test of finitely many linear inequalities; see [25, Section VII] for detailed discussions. Consequently, once the eigenvalues of A are found, not only is the verification problem finitely verifiable, it can also be solved by efficient linear programming methods.

3.4 Extension to Generalized Suplinear Sets

We briefly discuss an extension of positive invariant analysis to a more general set of the form $\mathcal{S} = \{x \in \mathbb{R}^n \mid H(x) \geq b\}$, where the vector-valued function H will be specified below. An extended real-valued function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is called *generalized sublinear* if (i) f is subadditive, i.e. $f(x + y) \leq f(x) + f(y), \forall x, y \in \mathbb{R}^n$, and (ii) f is generalized positively homogeneous, i.e., there exists a bijective and strictly increasing function $p : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ such that $f(\lambda x) = p(\lambda)f(x), \forall x \in \mathbb{R}^n, \forall \lambda \in \mathbb{R}_+$. It is clear that p is continuous on \mathbb{R}_+ . If $p(\lambda) \equiv \lambda$, then f becomes the standard sublinear function and is convex [10, 19]. Letting $\text{dom} f := \{x \in \mathbb{R}^n \mid f(x) \in \mathbb{R}\}$ be the domain of f , we call f *nontrivial* if $\text{dom} f \neq \emptyset$ and f is not identically zero on $\text{dom} f$. The following lemma states basic properties of p and f whose proof is given in the Appendix:

Lemma 22. Let the generalized sublinear function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be nontrivial. Then $f(0) = 0$ and p is inverse symmetric, i.e., $p(\frac{1}{\lambda}) = \frac{1}{p(\lambda)}, \forall \lambda > 0$ with $p(1) = 1$.

An example of a nontrivial generalized sublinear function is as follows with $p(\lambda) = \lambda^n$:

$$f(x) = \begin{cases} -\prod_{i=1}^n x_i, & \text{if } x \in \mathbb{R}_+^n \\ +\infty, & \text{otherwise} \end{cases}$$

A function $h : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{-\infty\}$ is called *generalized suplinear* if $-h$ is generalized sublinear. A vector-valued function $H : \mathbb{R}^n \rightarrow (\mathbb{R} \cup \{-\infty\})^m$ is called generalized suplinear if each of its components $H_j(x)$ is a generalized suplinear function defined by a (possibly different) bijective and strictly increasing function p_j . We assume that each $-H_j$ is nontrivial through this section. We call the set $\mathcal{S} \equiv \{x \in \mathbb{R}^n \mid H(x) \geq b\}$ *generalized suplinear* if H is generalized suplinear. Note that such a set is generally non-convex. Given a vector $u \in \mathbb{R}^n$, we write it as $u = \sum_i (0, \dots, 0, (u^i)^T, 0, \dots, 0)^T$ where $u^i \in \mathbb{R}^{n_i}$. Define $H^i(u^i) \equiv H((0, \dots, 0, (u^i)^T, 0, \dots, 0)^T)$. It is clear that $H^i : \mathbb{R}^{n_i} \rightarrow (\mathbb{R} \cup \{-\infty\})^m$ is generalized suplinear and $H(u) \geq \sum_i H^i(u^i)$ (due to subadditivity of $-H$). Moreover, $H^i(0) = 0$ which implies that its domain is nonempty. We also assume throughout this section that each component of H^i is not identically zero on its domain such that it is nontrivial.

The positive invariance results in the previous sections can be extended to yield sufficient existence conditions for generalized suplinear sets under the above assumptions. To illustrate this, we consider the linear dynamics $\dot{x} = Ax$ to avoid being excessively technical. Let $v^i \in V(\lambda_i)$ for a nonnegative eigenvalue λ_i of A such that $v^i = \bigoplus_j v^{ij}$. For an index set $\theta \subseteq \{1, \dots, m\}$, we have

$$\begin{aligned} H_\theta^i(e^{J_i t} v^i) &= H_\theta^i\left(e^{\lambda_i t} \sum_{k=0}^{\bar{n}_i-1} \frac{t^k}{k!} \bigoplus_j \mathcal{L}^k(v^{ij})\right) \geq \sum_{k=0}^{\bar{n}_i-1} H_\theta^i\left(e^{\lambda_i t} \frac{t^k}{k!} \bigoplus_j \mathcal{L}^k(v^{ij})\right) \\ &\geq \sum_{k=0}^{\bar{n}_i-1} \text{diag}(p_j(e^{\lambda_i t} \frac{t^k}{k!}))_{j \in \theta} H_\theta^i\left(\bigoplus_j \mathcal{L}^k(v^{ij})\right) \end{aligned}$$

where $\text{diag}(\rho_i)_{i \in \theta}$ denotes a diagonal matrix with diagonal entries given by indexed scalars $\rho_i, i \in \theta$. Define the tuple

$$P(H_\theta^i(v^i)) \equiv \left(H_\theta^i(\bigoplus_j \mathcal{L}^{\bar{n}_i-1}(v^{ij})), \dots, H_\theta^i(\bigoplus_j \mathcal{L}^0(v^{ij})) \right)$$

We call $P(H_\theta^i(v^i))$ *proper* if each entry in the tuple is finite. It can be shown that if a proper tuple $P(H_\theta^i(v^i)) \succ 0$ (resp. $P(H_\theta^i(v^i)) \equiv 0$), then there exists $u^i \in V(\lambda_i)$ such that $H_\theta^i(e^{J_i t} u^i) > 0$ (resp. $H_\theta^i(e^{J_i t} u^i) \geq 0$) for all $t \geq 0$ and any $a \in \mathbb{R}_{++}^m$. Indeed, similar to Proposition 8, let $u^i = \sum_{s=0}^{\bar{n}_i-1} \mu_s \bigoplus_j v^{ij}$ with nonnegative reals μ_s to be determined. If $P(H_\theta^i(v^i))$ is proper and $P(H_\theta^i(v^i)) \succ 0$, then

$$\begin{aligned} H_\theta^i(e^{J_i t} u^i) &\geq \sum_{k=0}^{\bar{n}_i-1} H_\theta^i\left(e^{\lambda_i t} \frac{t^k}{k!} \sum_{s=0}^{\bar{n}_i-1-k} \mu_s \bigoplus_j \mathcal{L}^{k+s}(v^{ij})\right) \\ &\geq \sum_{k=0}^{\bar{n}_i-1} \text{diag}(p_j(e^{\lambda_i t} \frac{t^k}{k!}))_{j \in \theta} \sum_{p=k}^{\bar{n}_i-1} H_\theta^i\left(\mu_{p-k} \bigoplus_j \mathcal{L}^p(v^{ij})\right) \geq a, \quad \forall t \geq 0 \end{aligned}$$

for suitable $\mu_s \geq 0$ by Corollary 7 and Lemma 22. The case $P(H_\theta^i(v^i)) \equiv 0$ follows from the above argument and subadditivity of $-H_\theta^i$. Furthermore, if there exist real eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_\ell \geq 0$ and $v^i \in V(\lambda_i), i = 1, \dots, \ell$ such that $P(H_\theta^i(v^i))$ are all proper and $P(H_\theta^1(v^1)) \succ 0$ for some index set θ , then for any $a \in \mathbb{R}_{++}^m$ and any given $u^i = \sum_{s=0}^{\bar{n}_i-1} \mu_{i s} \bigoplus_j v^{ij}, i \geq 2$ with $\mu_{i s} \geq 0$, we

obtain $u^1 = \sum_{s=0}^{\bar{n}_1-1} \mu_{1s} \bigoplus_j v^{1j} \in V(\lambda_1)$ with suitable $\mu_{1s} \geq 0$ such that

$$\begin{aligned} H_\theta(\bigoplus_{i=1}^{\ell} e^{J_i t} u^i) &\geq \sum_{i=1}^{\ell} H_\theta^i(e^{J_i t} u^i) \geq \text{diag}(p_j(e^{\lambda_1 t}))_{j \in \theta} \left[H_\theta^1 \left(\sum_{k=0}^{\bar{n}_1-1} \frac{t^k}{k!} \bigoplus_j \mathcal{L}^k(u^{1j}) \right) \right. \\ &\quad \left. + \sum_{i=2}^{\ell} \sum_{k=0}^{\bar{n}_i-1} \text{diag}(p_j(e^{(\lambda_i - \lambda_1)t} \frac{t^k}{k!}))_{j \in \theta} H_\theta^i \left(\bigoplus_j \mathcal{L}^k(u^{ij}) \right) \right] \geq a, \quad \forall t \geq 0 \end{aligned}$$

where we use the inverse symmetry of p_j 's shown in Lemma 22 and boundedness of $p_j(e^{(\lambda_i - \lambda_1)t} \frac{t^k}{k!})$ on $[0, \infty)$. Based on these results, we obtain the following sufficient conditions via the similar argument in (S1) of Theorem 12:

Proposition 23. Consider a nonempty generalized suplinear set \mathcal{S} with $b \geq 0$, i.e., $\gamma = \emptyset$. Then the positively invariant set \mathcal{A} of the linear dynamics on \mathcal{S} is nonempty if there exist real eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_k \geq 0$ and $v^i \in V(\lambda_i)$, $i = 1, \dots, k$ such that $P(H_\alpha^i(v^i))$ and $P(H_\beta^i(v^i))$ are proper for each i , $(P(H_\alpha^1(v^1)), \dots, P(H_\alpha^k(v^k))) \succ 0$, and $(P(H_\beta^1(v^1)), \dots, P(H_\beta^k(v^k))) \succneq 0$.

The sufficient conditions for the other cases can be established in the similar manner. Instead of further pursuing this generalization, we show two applications of the positive invariance results in the next two sections.

4 Application to Global Switching Characterization of Piecewise Affine Systems

In this section, the positive invariance results are applied to a class of affine hybrid systems. Specifically, necessary and sufficient conditions are derived to characterize global long-time switching behaviors of piecewise affine systems with isolated equilibria and infinite mode switchings.

4.1 Piecewise Affine Systems

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is called piecewise affine (PA) if there exists a finite family of affine functions $\{f_i\}_{i=1}^{\ell}$ such that $f(x) \in \{f_i(x)\}_{i=1}^{\ell}$ for each $x \in \mathbb{R}^n$ [15, 31]. Consider the ODE system: $\dot{x} = f(x)$, where $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is continuous and piecewise affine. We call such a system the *piecewise affine system* (PAS) [11]. A continuous and PA function possesses an appealing geometric structure for its domain, which provides an alternative representation for the PAS. To elaborate on this, let Ξ be a finite collection of polyhedra $\mathcal{P}_i = \{x \in \mathbb{R}^n \mid C_i x \geq b_i\}$ for $C_i \in \mathbb{R}^{m_i \times n}$ and $b_i \in \mathbb{R}^{m_i}$. A face of \mathcal{P}_i is defined as $\mathcal{P}_i \cap \{x \mid (C_i x - b_i)_\alpha = 0\}$ for a nonempty index set α such that there exists $x \in \mathbb{R}^n$ with $(C_i x - b_i)_\alpha = 0$ and $(C_i x - b_i)_{\bar{\alpha}} > 0$, where $\bar{\alpha}$ denotes the complement of α ; see [31, Proposition 2.1.3] for more details. A face of \mathcal{P}_i is called proper if it does not coincide with \mathcal{P}_i . We call Ξ a *polyhedral subdivision* of \mathbb{R}^n [15, 31] if

- (a) the union of all polyhedra in Ξ is equal to \mathbb{R}^n , i.e., $\bigcup_{i=1}^m \mathcal{P}_i = \mathbb{R}^n$,
- (b) each polyhedron in Ξ has a nonempty interior (thus is of dimension n), and
- (c) the intersection of any two polyhedra in Ξ is either empty or a common proper face of both polyhedra, i.e., $\mathcal{P}_i \cap \mathcal{P}_j \neq \emptyset \implies [\mathcal{P}_i \cap \mathcal{P}_j = \mathcal{P}_i \cap \{x \mid (C_i x - b_i)_\alpha = 0\} = \mathcal{P}_j \cap \{x \mid (C_j x - b_j)_\beta = 0\}]$ for nonempty index sets α and β with $\mathcal{P}_i \cap \{x \mid (C_i x - b_i)_\alpha = 0\} \neq \mathcal{P}_i$ and $\mathcal{P}_j \cap \{x \mid (C_j x - b_j)_\beta = 0\} \neq \mathcal{P}_j$.

For a continuous and PA function f , one can always find a polyhedral subdivision of \mathbb{R}^n and finitely many affine functions $g_i \equiv A_i x + d_i$ such that f coincides with one of g_i 's on each polyhedron in Ξ [15, Proposition 4.2.1]. With these notions, we can rewrite the PAS in the equivalent form

$$\dot{x} = A_i x + d_i \quad \text{if } x \in \mathcal{P}_i, \quad (9)$$

where $[x \in \mathcal{P}_i \cap \mathcal{P}_j] \Rightarrow [A_i x + d_i = A_j x + d_j]$ holds true due to the continuity of f . In what follows, we call each affine dynamics $\dot{x} = A_i x + d_i$ and its associated polyhedron \mathcal{P}_i a *mode* of the PAS. Since a continuous and PA function is globally Lipschitz continuous [15], the PAS (9) has a unique solution for any initial state. The PASs form a class of affine hybrid systems, for which the time-invariant vector fields are affine, the invariant sets are the polyhedrons \mathcal{P}_i of dimension n , the guard sets are the boundaries of these polyhedra, and the reset maps are all identities. The PASs model important hybrid systems in various areas and thus have attracted considerable research attention. Examples include affine complementarity systems in dynamical optimization [26, 32] and nonsmooth dynamical systems such as genetic regulatory networks in systems biology [14].

An important class of the PASs is the conewise linear systems (CLSs) [11, 34, 35], for which $b_i = 0$ and $d_i = 0$ for each i such that (9) becomes $\dot{x} = A_i x$ if $x \in \mathcal{C}_i \equiv \{x \mid C_i x \geq 0\}$. We shall discuss more about the CLSs and PASs in the next subsection.

4.2 Mode Switchings and Global Long-time Switching Behaviors

A PAS is subject to state-dependent mode switchings along its trajectories with implicit transition times and implicit mode selection at switching times. The mode switching is formally defined as

Definition 24. For a given state trajectory $x(t, x^0)$ of the PAS (9), we say that a time instant $t_* \geq 0$ is *not* a switching time along $x(t, x^0)$ if there exist $i \in \{1, \dots, m\}$ and $\varepsilon > 0$ such that $x(t, x^0) \in \mathcal{P}_i, \forall t \in [t_* - \varepsilon, t_* + \varepsilon]$; otherwise, t_* is a *switching time* along $x(t, x^0)$, and that the PAS has a *mode transition* or *mode switching* at t_* along $x(t, x^0)$.

While state-dependent mode transitions are fairly common in applications with examples including nonsmooth physical systems and dynamic optimization, they usually yield complicate dynamical and control issues for hybrid dynamics. As a special case of the PASs, the CLSs have been extensively studied in the recent papers [11, 34], particularly for their state-dependent mode switching behaviors. In specific, a CLS is shown to be Zeno free [11] and possess the ‘‘simple switching property’’ [34]. These results are critical to various local and global dynamic properties of the CLSs. We extend these results to the PASs in this section, motivated by explicit characterization of long-time switching behaviors of the PASs. It can be shown that the PASs enjoy similar local switching features as the CLSs, e.g. the non-Zenoness and simple switching property; see the Appendix for their formal definitions and proofs as well as other background results. However, the extension to long-time switching properties turns out to be rather nontrivial since these properties are closely related to positively invariant sets of affine modes of the PASs that have not been fully treated. The positive invariance results in Section 3 provide essential tools that lead to verifiable characterization conditions as shown below.

Let \mathcal{A}_i and $\mathcal{E}_i \equiv \{x \mid A_i x + d_i = 0, x \in \mathcal{P}_i\}$ denote the positively invariant set and equilibrium set of the i th mode of the PAS (9), respectively. An equilibrium x^e of the PAS is called *isolated* if a neighborhood of x^e exists such that it does not contain any other equilibrium. We call a PAS *with isolated equilibria* if each of its equilibria is isolated. Notice that an equilibrium is asymptotically stable only if it is isolated. Hence a PAS with isolated equilibria is of special interest in asymptotic stability analysis. It is easy to show via the convex structure of the PAS that each equilibrium of the PAS is isolated if and only if the equilibrium set of each mode contains at most one element. Therefore, such a PAS has finitely many equilibria.

In the sequel, we concentrate on a specific class of PASs, namely, the PASs whose trajectories have infinitely many mode transitions (in the positive time direction) whenever they start from non-equilibria. These PASs are referred to as the PASs with infinite mode switchings. The following result characterizes such a PAS via the positively invariant set of each mode. Its proof, relying on the non-Zeno and simple switching properties, is given in the Appendix.

Lemma 25. The PAS (9) has infinite mode switchings if and only if $\mathcal{A}_i = \mathcal{E}_i$ for each i .

In view of this lemma, we focus on each affine mode of the PAS. For the i th mode and its equilibrium set \mathcal{E}_i , we assume the (possibly zero) d_i to be in the null space of A_i^T as before. Obviously if \mathcal{E}_i is nonempty, then d_i must be zero. Consider the following cases:

(1) $\mathcal{E}_i = \emptyset$. In this case, $\mathcal{A}_i = \mathcal{E}_i$ if and only if exactly one of the following holds: (1.1) the existence conditions in Theorem 12 fail when $d_i = 0$; (1.2) the existence conditions in Theorem 18 do not hold when $d_i \neq 0$.

(2) \mathcal{E}_i is a singleton set. In this case, since $d_i = 0$, the sole equilibrium in \mathcal{E}_i can be taken as $x^e = 0$ (by a suitable state translation). Therefore, $b_i \leq C_i x^e = 0$, which implies that b_i has no positive element or equivalently the index set α associated with b_i is empty. Hence, $\mathcal{A}_i = \mathcal{E}_i$ if and only if \mathcal{A}_i is trivial, i.e., $\mathcal{A}_i = \{0\}$, or equivalently the nontrivial positive invariance conditions in Theorem 16 fail.

By virtue of the above discussions, we immediately obtain the following result without proof:

Corollary 26. A PAS with isolated equilibria has infinite mode switchings if and only if the conditions in (1) or (2) discussed above hold for each of its modes.

Remark 27. We make a few comments before closing this section: (i) For a CLS, all of its equilibria are isolated if and only if it has a unique equilibrium at the origin. Hence, we deduce from Theorem 16 that a CLS with isolated equilibria has infinite mode switchings if and only if \mathcal{A}_i is trivial for each i , which is further equivalent to the conditions that (C_i, A_i) is an observable pair and that A_i has no eigenvector in the polyhedral cone $\{x \mid C_i x \geq 0\}$. This recovers the result of [34, Proposition 17]. (ii) if \mathcal{E}_i is non-degenerate, i.e., $\mathcal{E}_i \cap \text{int } \mathcal{P}_i$ is nonempty, then the similar characterization conditions for $\mathcal{A}_i = \mathcal{E}_i$ can be obtained as those in [34, Proposition 19] via the positive invariant results and the observation that $d_i = 0$.

5 Application to Exact Safety Verification of Affine Dynamics

A general safety verification problem poses a challenging analytical and numerical problem, even for relatively simpler dynamics and constraint sets, because of infinite dimensional nature of ODE dynamics. For this reason, two technical paths have been widely followed in the literature: one is based on approximation methods (e.g., over-approximation/under-approximation and asymptotic approximation) for general nonlinear dynamics or fast computation [28, 29, 39], and the other focuses on exact approaches but only for simpler dynamics such as linear or affine dynamics [22, 23, 40]. We take the second path in the current paper and concentrate on exact safety verification of affine dynamics. The affine dynamics and algebraic structure of constraint sets allow us to obtain less conservative and computationally tractable verification results. This is demonstrated by the following example.

Example 28. Consider the affine dynamics $\dot{x} = Ax + d$ with the polyhedral initial set \mathcal{S}_0 and the final safe region \mathcal{S}_f . By Minkowski–Weyl Decomposition Theorem, we decompose \mathcal{S}_0 into the sum of a compact convex hull and a conic hull, i.e., $\mathcal{S}_0 = \text{conv}(v^1, \dots, v^\ell) + \text{cone}(u^1, \dots, u^k)$ for the extreme points v^i and the extreme rays u^j , where conv and cone denote the convex hull and the closed conic hull of the given sets, respectively. Using this decomposition and letting

$\mathcal{S}_f = \{x \mid Cx \geq b\}$, it is easy to verify that the affine dynamics is safe from \mathcal{S}_0 on $\Delta = \mathbb{R}_+$ if and only if $C[e^{At} v^i + \int_0^t e^{A(t-\tau)} d\tau d] \geq b, \forall t \geq 0$ for all $i = 1, \dots, \ell$ and $C[e^{At} w^j + \int_0^t e^{A(t-\tau)} d\tau d] \geq 0, \forall t \geq 0$ for all $j = 1, \dots, k$. These conditions enable one to check for finitely many vectors only and thus considerably simplify computations.

Analytic results are hard to obtain for general \mathcal{S}_0 and \mathcal{S}_f in safety verification, and the best way to solve the (exact) safety verification problem is to pursue numerical approaches. This yields a critical computability question, namely, whether such the problem is finitely verifiable or decidable [8]. A typical approach to address this issue is to convert the original safety verification problem into a semi-algebraic problem and then apply the celebrated Tarski-Seidenberg decision procedure [9, Corollary 1.4.7]. Such an approach has been successfully invoked to show decidability of several classes of linear dynamics on semi-algebraic sets, for example, [13, 22, 23, 40].

In the following, we show that the dynamic analysis techniques for positive invariance analysis, along with semi-algebraic arguments, lead to improved decidability results. Specifically, we consider an affine dynamics whose defining matrix contains complex eigenvalues only. Examples of such a system include linear Hamiltonian dynamics with complex eigenvalues only [3, Appendix 6]. Furthermore, we assume that \mathcal{S}_0 and \mathcal{S}_f are closed semi-algebraic sets, i.e., they are described by finitely many multi-variate polynomial equations or inequalities which are neither convex nor polyhedral in general. It is shown below that the safety verification problem is decidable even if the ratios of mode frequencies of the dynamics are irrational. Two key tools have been exploited to establish this result: (i) Lemmas 3 and 5 that express the infimum of a linear combination of periodic functions in term of the sum of the minima of the periodic functions, and (ii) an algebraic technique that formulates finding the minima of the periodic functions as a semi-algebraic problem.

Admittedly, the obtained decidability result is perhaps mainly of theoretical interest since practical computations assume rational numbers. Nevertheless it reveals an interesting perspective for decidability analysis via dynamic system techniques that may be extended to a broader class of systems.

Proposition 29. Let $\mathcal{S}_0 = \{x \in \mathbb{R}^n \mid p(x) \geq 0, w(x) = 0\}$ and $\mathcal{S}_f = \{x \in \mathbb{R}^n \mid f(x) \geq 0\}$ be closed semi-algebraic sets, where $p, w,$ and f are vector-valued multivariate polynomials. Suppose that A has only complex eigenvalues which are all known. Then checking safety of the affine dynamics $\dot{x} = Ax + d$ on the time interval $\Delta = \mathbb{R}$ is decidable.

Proof. Since A has no real eigenvalue, it is invertible and thus d is always in the range of A . Therefore, via a suitable affine coordinate transformation, we assume d to be zero such that the affine dynamics becomes linear, i.e., $\dot{x} = Ax$. The sets \mathcal{S}_0 and \mathcal{S}_f remain semialgebraic after the transformation. Since all the complex eigenvalues of A are known, each entry of $e^{At}x$ can be written as $e^{\lambda t} t^p l(x) \cos(\omega t)$ or $e^{\lambda t} t^p l(x) \sin(\omega t)$ for the known λ, p, ω , where $p \geq 0$ is an integer, $\omega > 0$, $\lambda \pm i\omega$ is the corresponding complex eigenvalue of A , and $l(x)$ is a linear function of x . Therefore, via basic trigonometric relations and straightforward computations, we have $f(e^{At}x) = c(x) + \sum_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})} e^{\tilde{\lambda} t} t^{\tilde{p}} [g_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x) \cos(\tilde{\omega} t) + h_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x) \sin(\tilde{\omega} t)]$, where $c(x)$ is a (vector-valued) multivariate polynomial, $g_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x), h_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x)$ are (vector-valued) multivariate polynomials corresponding to the tuple $(\tilde{\lambda}, \tilde{p}, \tilde{\omega})$, and $\tilde{\omega}$'s associated with the same $(\tilde{\lambda}, \tilde{p})$ are all distinct. Notice that there are finitely many such tuples. We have the following claim:

Claim C1: the linear system is safe on $\Delta = \mathbb{R}$ if and only if for all $x \in \mathcal{S}_0$, $g_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x) = h_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x) = 0$ for $(\tilde{\lambda}, \tilde{p}) \neq 0$ and $c(x) + \sum_{\omega_j} [g_{(0,0,\omega_j)}(x) \cos(\omega_j t) + h_{(0,0,\omega_j)}(x) \sin(\omega_j t)] \geq 0, \forall t \in \mathbb{R}$, where each ω_j in the latter condition corresponds to the pure imaginary eigenvalues of A .

The proof for (C1) is as follows. The sufficiency is obvious. To see necessity, we first show that for each $x \in \mathcal{S}_0$, $g_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x) = h_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x) = 0$ whenever $(\tilde{\lambda}, \tilde{p}) \succ (0, 0)$. Suppose not. Then there

exist $x^* \in \mathcal{S}_0$ and $(\tilde{\lambda}, \tilde{p}) \succ (0, 0)$ with $(g_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x^*), h_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x^*)) \neq 0$ for some $\tilde{\omega} > 0$. Let (λ^*, p^*) be the largest such a pair in the lexicographical sense, i.e., $(\lambda^*, p^*) \succ (\tilde{\lambda}, \tilde{p})$ for any pair $(\tilde{\lambda}, \tilde{p})$ with nonzero $(g_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x^*), h_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x^*))$. Therefore, $f(e^{At}x^*)$ tends to $e^{\lambda^* t} t^{p^*} \sum_{\tilde{\omega}} [g_{(\lambda^*, p^*, \tilde{\omega})}(x^*) \cos(\tilde{\omega}t) + h_{(\lambda^*, p^*, \tilde{\omega})}(x^*) \sin(\tilde{\omega}t)]$ as $t \rightarrow +\infty$. However, by Corollary 4 (or Proposition 11) and $f(e^{At}x^*) \geq 0, \forall t \geq 0$, we have $g_{(\lambda^*, p^*, \tilde{\omega})}(x^*) = h_{(\lambda^*, p^*, \tilde{\omega})}(x^*) = 0$ for all $\tilde{\omega}$, a contradiction. Similarly, by the reverse-time argument, we conclude that if $(\tilde{\lambda}, \tilde{p}) \prec (0, 0)$ (or equivalently $\tilde{\lambda} < 0$), $g_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x) = h_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x) = 0$ for all $x \in \mathcal{S}_0$. Hence, the claim holds.

In the sequel, we show that the necessary and sufficient conditions established by Claim C1 can be finitely verified. It is easy to see that checking the first condition, i.e., $g_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x) = h_{(\tilde{\lambda}, \tilde{p}, \tilde{\omega})}(x) = 0$ with $(\tilde{\lambda}, \tilde{p}) \neq 0$ for all $x \in \mathcal{S}_0$, can be cast as a semialgebraic decision problem and thus is decidable. Hence, we only focus on the second condition which possesses a universal quantifier t . For the ease of development, we consider two cases as follows:

(1) The ratio of any two frequencies associated with pure imaginary eigenvalues of A is rational. This case follows from [40, Section IV].

(2) Among all the frequencies ω_i of the pure imaginary eigenvalues of A , some of $\omega_i/\omega_j, i \neq j$, are rational and the others are irrational. Define the function $s_{\omega_j}(t, x) \equiv g_{(0,0,\omega_j)}(x) \cos(\omega_j t) + h_{(0,0,\omega_j)}(x) \sin(\omega_j t)$ for each ω_j . Following the similar treatment as in Section 2.2, we obtain the collection of (disjoint) equivalent classes $E_{\omega_j} = \{s_{\omega_j}(t, x) \mid \omega_i/\omega_j \text{ is rational}\}$ (which is independent of x). Let $\tilde{\omega}_\ell > 0$ be a basis frequency associated with each E_{ω_j} and denote such the equivalent class by $E_{\tilde{\omega}_\ell}$. Let $q_{\tilde{\omega}_\ell}(t, x) \equiv \sum_{s_{\omega_i} \in E_{\tilde{\omega}_\ell}} s_{\omega_i}(t, x)$. For any fixed $x \in \mathbb{R}^n$, we see that each component

of $q_{\tilde{\omega}_\ell}(t, x)$ enjoys the same four properties of $q_{\tilde{\omega}_s}$ stated in Section 2.2. In particular, if the i th component $(q_{\tilde{\omega}_\ell})_i(\cdot, x)$ is not identically zero, then it attains the maximal and minimal values $\mu_{i, \tilde{\omega}_\ell}(x) > 0$ and $\nu_{i, \tilde{\omega}_\ell}(x) < 0$ on $(-\infty, \infty)$ respectively. Furthermore, recall that for each $s_{\omega_j} \in E_{\tilde{\omega}_\ell}$, $\omega_j/\tilde{\omega}_\ell$ is a positive integer. Therefore, by the basic trigonometric results, we see that each $q_{\tilde{\omega}_\ell}(t, x)$ can be expressed as a (vector-valued) multivariate polynomial function in terms of x , $\sin(\tilde{\omega}_\ell t)$, and $\cos(\tilde{\omega}_\ell t)$. In other words, letting $u_{\tilde{\omega}_\ell} \equiv \sin(\tilde{\omega}_\ell t)$ and $v_{\tilde{\omega}_\ell} \equiv \cos(\tilde{\omega}_\ell t)$ for each $E_{\tilde{\omega}_\ell}$, we obtain a (vector-valued) polynomial function $\tilde{q}_{\tilde{\omega}_\ell}(x, u_{\tilde{\omega}_\ell}, v_{\tilde{\omega}_\ell})$ that is equivalent to $q_{\tilde{\omega}_\ell}(t, x)$, where $u_{\tilde{\omega}_\ell}^2 + v_{\tilde{\omega}_\ell}^2 = 1$. Likewise, we can write the time derivative $\frac{\partial q_{\tilde{\omega}_\ell}(t, x)}{\partial t}$ as an equivalent (vector-valued) polynomial function $\tilde{d}_{\tilde{\omega}_\ell}(x, u_{\tilde{\omega}_\ell}, v_{\tilde{\omega}_\ell})$. Suppose there are k equivalent classes $E_{\tilde{\omega}_\ell}$. Since $c(x) + \sum_{\omega_j} [g_{(0,0,\omega_j)}(x) \cos(\omega_j t) + h_{(0,0,\omega_j)}(x) \sin(\omega_j t)] = c(x) + \sum_{\ell=1}^k q_{\tilde{\omega}_\ell}(t, x)$, we have:

Claim C2: for any fixed x and each i , $c_i(x) + \sum_{\ell=1}^k (q_{\tilde{\omega}_\ell})_i(t, x) \geq 0, \forall t \in \mathbb{R}$ if and only if the following implication holds

$$\left[(\tilde{d}_{\tilde{\omega}_\ell})_i(x, u_{\tilde{\omega}_\ell}, v_{\tilde{\omega}_\ell}) = 0, u_{\tilde{\omega}_\ell}^2 + v_{\tilde{\omega}_\ell}^2 = 1, \forall \ell = 1, \dots, k \right] \Rightarrow c_i(x) + \sum_{\ell=1}^k (\tilde{q}_{\tilde{\omega}_\ell})_i(x, u_{\tilde{\omega}_\ell}, v_{\tilde{\omega}_\ell}) \geq 0 \quad (10)$$

The proof for Claim C2 is as follows:

“Sufficiency”: Without loss of generality, we assume that $(q_{\tilde{\omega}_\ell})_i(\cdot, x)$ is not identically zero for each ℓ . Observe that for every ℓ , the real pairs $(u_{\tilde{\omega}_\ell}, v_{\tilde{\omega}_\ell})$ that satisfy $(\tilde{d}_{\tilde{\omega}_\ell})_i(x, u_{\tilde{\omega}_\ell}, v_{\tilde{\omega}_\ell}) = 0$ and $u_{\tilde{\omega}_\ell}^2 + v_{\tilde{\omega}_\ell}^2 = 1$ correspond to critical points of the real-valued function $(q_{\tilde{\omega}_\ell})_i(\cdot, x)$ (by the definition of $\tilde{d}_{\tilde{\omega}_\ell}$). Therefore, one of such the $(u_{\tilde{\omega}_\ell}, v_{\tilde{\omega}_\ell})$'s, say $(u_{\tilde{\omega}_\ell}^*, v_{\tilde{\omega}_\ell}^*)$, is a minimizer of the bounded periodic function $(q_{\tilde{\omega}_\ell})_i(\cdot, x)$. This shows that $(\tilde{q}_{\tilde{\omega}_\ell})_i(x, u_{\tilde{\omega}_\ell}^*, v_{\tilde{\omega}_\ell}^*) = \nu_{i, \tilde{\omega}_\ell}(x) < 0$ (see the definition of $\nu_{i, \tilde{\omega}_\ell}$ above). We thus deduce from (10) that $c_i(x) + \sum_{\ell=1}^k \nu_{i, \tilde{\omega}_\ell}(x) \geq 0$. Since $(q_{\tilde{\omega}_\ell})_i(t, x) \geq \nu_{i, \tilde{\omega}_\ell}(x)$

for all t , we have $c_i(x) + \sum_{\ell=1}^k (q_{\tilde{\omega}_\ell})_i(t, x) \geq c_i(x) + \sum_{\ell=1}^k \nu_{i, \tilde{\omega}_\ell}(x) \geq 0$ for all t as desired.

“Necessity”: We show this by contradiction. Suppose $c_i(x) + \sum_{\ell=1}^k q_{\tilde{\omega}_\ell}(t, x) \geq 0, \forall t \in \mathbb{R}$ but there exist pairs $(u_{\tilde{\omega}_\ell}^*, v_{\tilde{\omega}_\ell}^*)$ with $(u_{\tilde{\omega}_\ell}^*)^2 + (v_{\tilde{\omega}_\ell}^*)^2 = 1$ and $(\tilde{d}_{\tilde{\omega}_\ell})_i(x, u_{\tilde{\omega}_\ell}^*, v_{\tilde{\omega}_\ell}^*) = 0, \ell = 1, \dots, k$, such that $c_i(x) + \sum_{\ell=1}^k (\tilde{q}_{\tilde{\omega}_\ell})_i(x, u_{\tilde{\omega}_\ell}^*, v_{\tilde{\omega}_\ell}^*) < 0$. Notice that for each $\ell, (\tilde{q}_{\tilde{\omega}_\ell})_i(x, u_{\tilde{\omega}_\ell}^*, v_{\tilde{\omega}_\ell}^*) \in [\nu_{i, \tilde{\omega}_\ell}(x), \mu_{i, \omega_\ell}(x)]$ and that $(q_{\tilde{\omega}_\ell})_i(t, x)$ is onto $[\nu_{i, \tilde{\omega}_\ell}(x), \mu_{i, \tilde{\omega}_\ell}(x)]$. Consequently, by the properties of $(q_{\tilde{\omega}_\ell})_i(\cdot, x)$ and Lemmas 3 and 5 (which rely on irrational ratio of any two basis frequencies), we deduce the existence of $t_* \in \mathbb{R}$ such that $(q_{\tilde{\omega}_\ell})_i(t_*, x)$ is arbitrarily close to $(\tilde{q}_{\tilde{\omega}_\ell})_i(x, u_{\tilde{\omega}_\ell}^*, v_{\tilde{\omega}_\ell}^*)$ for each $\ell = 1, \dots, k$. This, together with the continuity of $(q_{\tilde{\omega}_\ell})_i(\cdot, x)$, shows that $c_i(x) + \sum_{\ell=1}^k (q_{\tilde{\omega}_\ell})_i(t_*, x) < 0$, which is a contradiction.

In view of Claim C2, we now complete the proof for the second case. Recalling $x \in \mathcal{S}_0 \Leftrightarrow [p(x) \geq 0, w(x) = 0]$ and defining $\hat{u}^i, \hat{v}^i \in \mathbb{R}^k, i = 1, \dots, m$, we obtain the following implication that is equivalent to the second condition of Claim C1: for each i ,

$$[p(x) \geq 0, w(x) = 0, (\tilde{d}_{\tilde{\omega}_\ell})_i(x, \hat{u}_j^i, \hat{v}_j^i) = 0, (\hat{u}_j^i)^2 + (\hat{v}_j^i)^2 - 1 = 0] \implies c_i(x) + \sum_{\ell=1}^k (\tilde{q}_{\tilde{\omega}_\ell})_i(x, \hat{u}_j^i, \hat{v}_j^i) \geq 0$$

Denoting $\tilde{x} = (x, \hat{u}^1, \dots, \hat{u}^m, \hat{v}^1, \dots, \hat{v}^m) \in \mathbb{R}^{n+k \times 2m}$, we can rewrite the above implication as

$$[\tilde{p}(\tilde{x}) \geq 0, \tilde{w}(\tilde{x}) = 0] \implies [\tilde{g}(\tilde{x}) = 0, \tilde{h}(\tilde{x}) \geq 0],$$

where $\tilde{p}, \tilde{w}, \tilde{g}$ and \tilde{h} are appropriate vector-valued polynomials. Since checking this implication can be accomplished in finite steps via the Tarski-Seidenberg decision procedure, the safety verification problem is decidable. \square

The above result can be easily extended to a broader semi-algebraic set of the form $\mathcal{S}_f = \{x \in \mathbb{R}^n \mid f(x) \geq 0, q(x) = 0\}$ with polynomials f and q , since \mathcal{S}_f is equivalent to $\{x \in \mathbb{R}^n \mid f(x) \geq 0, q(x) \geq 0, -q(x) \geq 0\}$. The following illustrative example shows how to convert a safety verification problem into a decidable semi-algebraic problem via Proposition 29.

Example 30. Consider the linear system on \mathbb{R}^8 whose defining matrix $A = \text{diag}(A_1, A_2, A_3, A_4)$. Here the matrix blocks A_i are

$$A_1 = \begin{bmatrix} \sigma_1 & \omega_1 \\ -\omega_1 & \sigma_1 \end{bmatrix}, \quad A_2 = \begin{bmatrix} 0 & \pi \\ -\pi & 0 \end{bmatrix}, \quad A_3 = \begin{bmatrix} 0 & 2\pi \\ -2\pi & 0 \end{bmatrix}, \quad A_4 = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix},$$

where $\sigma_1 \neq 0$ and $\omega_1 > 0$. Let the initial set $\mathcal{S}_0 = \{x \in \mathbb{R}^8 \mid \|x - x^*\|_2 \leq 1\}$ for a given x^* and the safe region $\mathcal{S}_f = \{x \in \mathbb{R}^8 \mid c^T x \geq b\}$ for some $c \in \mathbb{R}^8$ and $b \in \mathbb{R}$. To simplify notation, let $c^T = (c_1^T, \dots, c_4^T)$ and $x = ((x^1)^T, \dots, (x^4)^T)^T$, where $c_i, x^i \in \mathbb{R}^2$ correspond to the matrix block A_i , and let the symplectic matrix $S = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$. Therefore,

$$c^T e^{At} x = e^{\sigma_1 t} [c_1^T x^1 \cos(\omega_1 t) + (S c_1)^T x^1 \sin(\omega_1 t)] + q_1(x, t) + q_2(x, t),$$

where $q_1(t, x) = c_2^T x^2 \cos(\pi t) + (S c_2)^T x^2 \sin(\pi t) + c_3^T x^3 \cos(2\pi t) + (S c_3)^T x^3 \sin(2\pi t)$ and $q_2(t, x) = c_4^T x^4 \cos(t) + (S c_4)^T x^4 \sin(t)$. Notice that $q_1(t, x)$ and $q_2(t, x)$ are periodic in t but their frequency ratio is irrational. Moreover, for a fixed x , even though q_1 is the sum of two sinusoidal functions with frequencies π and 2π respectively, the maximal (resp. minimal) values of q_1 cannot be simply written as the sum of the maximal (resp. minimal) values of the two sinusoidal functions. This is why we introduce the time derivative of q_1 to characterize the extremal value of q_1 . Now define

$u_1 \equiv \cos(\pi t), v_1 \equiv \sin(\pi t)$ and $u_2 \equiv \cos(t), v_2 \equiv \sin(t)$. Rewriting $q_i(t, x)$ and $\frac{\partial q_i(t, x)}{\partial t}, i = 1, 2$ in terms of x, u_1, v_1, u_2, v_2 and applying the argument in Proposition 29, we transform the safety verification problem on the time domain $\Delta = \mathbb{R}$ into the following semi-algebraic decision problem $[P \implies Q]$ on \mathbb{R}^{12} , where

$$P \equiv [\pi(-c_2^T x^2 v_1 + (Sc_2)^T x^2 u_1 - 4c_3^T x^3 u_1 v_1 + 2(Sc_3)^T x^3 (u_1^2 - v_1^2)) = 0] \wedge [u_1^2 + v_1^2 - 1 = 0] \\ \wedge [-c_4^T x^4 v_2 + (Sc_4)^T x^4 u_2 = 0] \wedge [u_2^2 + v_2^2 - 1 = 0] \wedge [(x - x^*)^T (x - x^*) - 1 \leq 0],$$

and

$$Q \equiv [c_2^T x^2 u_1 + (Sc_2)^T x^2 v_1 + c_3^T x^3 (u_1^2 - v_1^2) + 2(Sc_3)^T x^3 u_1 v_1 + c_4^T x^4 u_2 + (Sc_4)^T x^4 v_2 - b \geq 0] \\ \wedge [c_1^T x^1 = 0] \wedge [(Sc_1)^T x^1 = 0].$$

The latter problem is decidable and can be solved using the quantifier elimination technique. The recent sum-of-squares relaxation approach provides a numerically efficient alternative via powerful semidefinite programming techniques. Nevertheless further exploration of numerical issues is beyond the scope of the current paper. We refer the reader to [27, 28, 40] for additional information.

6 Conclusions

In this paper, we have addressed the existence of a positively invariant set of an affine dynamics on a convex polyhedron. Necessary and sufficient existence conditions are established via algebraic properties of the lexicographical relation and long-time dynamic analysis techniques. These positive invariance results form a cornerstone for long-time switching characterization of piecewise affine systems and safety verification of affine dynamics on semi-algebraic sets. A deeper understanding of switching behaviors of affine hybrid dynamics and safety analysis problems warrants a further investigation of the positive invariance problem, e.g., further characterization of algebraic and geometric structures of the positively invariant sets. Partial results along this line have been obtained by the author to understand positively invariant cones for the CLSs. Another interesting problem is how to develop refined and less conservative stability or other long-time dynamic results using the switching properties characterized in this paper. Preliminary results have been reported in [33, 34] for the CLSs and an extension to the PASs is expected. Feedback control design that ensures (non-)existence of a positively invariant set in a convex polyhedron remains an open, yet interesting, issue to be investigated. Finally, numerical aspects of positive invariance analysis and safety verification, e.g., complexity analysis and algorithm design, pose practically important problems that will be addressed in the future.

7 Appendix

7.1 Proof of Theorem 12

Proof. The case where the index set α is empty is trivial since $0 \in \mathcal{A}$, which shows that \mathcal{A} is nonempty. Note that the matrix A needs not to have a real eigenvalue in this case. In the subsequent development we assume that α is nonempty.

“Sufficiency”. We consider two cases as follows:

(S1) $\alpha \neq \emptyset, \gamma = \emptyset$. In this case, there exist real eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_k \geq 0$ and $v^i \in V(\lambda_i)$ such that $(P(C_\alpha^1, v^1), \dots, P(C_\alpha^k, v^k)) \succ 0$ and $(P(C_\beta^1, v^1), \dots, P(C_\beta^k, v^k)) \not\prec 0$. Let $\alpha \subseteq \theta \subseteq \alpha \cup \beta$ be such that $\tilde{P} \equiv (P(C_\theta^1, v^1), \dots, P(C_\theta^k, v^k)) \succ 0$ and $(P(C_{\bar{\theta}}^1, v^1), \dots, P(C_{\bar{\theta}}^k, v^k)) \equiv 0$, where $\bar{\theta}$ is the complement of θ . Hence, each row of \tilde{P} contains a positive leading entry. Without

loss of generality, we further assume that \tilde{P} is of the echelon structure via suitable row switchings and that each $P(C_{\theta}^i, v^i)$ contains a positive leading entry in \tilde{P} . For each $P(C_{\theta}^i, v^i)$, define the index set

$$\theta_i \equiv \{j \in \theta \mid \text{the } j\text{th row of } P(C_{\theta}^i, v^i) \text{ contains a positive leading entry of } \tilde{P}\}.$$

It is clear that θ_i is nonempty, $\theta = \cup_{i=1}^k \theta_i$, $\theta_i \cap \theta_j = \emptyset$ for any $i \neq j$. Besides, due to the lexicographical relation and echelon structure of \tilde{P} , we have, for each i , $P(C_{\theta_i}^i, v^i) \succ 0$ and $P(C_{\theta_i}^j, v^j) \equiv 0$ for $1 \leq j < i$. The following diagram illustrates this structure which is critical to the rest of the proof:

$$\tilde{P} = \begin{pmatrix} P(C_{\theta_1}^1, v^1) & \star & \cdots & \cdots & \star \\ & P(C_{\theta_2}^2, v^2) & \star & \cdots & \star \\ & & \ddots & & \vdots \\ & & & P(C_{\theta_{k-1}}^{k-1}, v^{k-1}) & \star \\ & & & & P(C_{\theta_k}^k, v^k) \end{pmatrix}$$

where \star denotes the blocks without a leading entry in \tilde{P} . For each i , define the index set $\tilde{\theta}_i \equiv \cup_{j=i}^k \theta_j$. We claim that for every $i = k, \dots, 1$ and any $h \in \mathbb{R}_{++}^{|\tilde{\theta}_i|}$, there exists a vector $u \equiv \bigoplus_{j=i}^k u^j \in \bigoplus_{j=i}^k V(\lambda_j)$ such that $\sum_{j=i}^k C_{\theta_i}^j e^{J_j t} u^j \geq h$ for all $t \geq 0$. We prove the claim by induction on i as follows. First, consider $i = k$. In this case, the claim follows immediately from Proposition 8. Suppose that the claim holds for all $i = k, \dots, \ell + 1$, where $1 \leq \ell \leq k - 1$. Now consider $i = \ell$. By the induction hypothesis, for a given $h = (h_1, h_2) \in \mathbb{R}_{++}^{|\tilde{\theta}_\ell|}$ with $h_1 \in \mathbb{R}_{++}^{|\theta_\ell|}$ and $h_2 \in \mathbb{R}_{++}^{|\tilde{\theta}_{\ell+1}|}$, there exists $\tilde{u} \equiv \bigoplus_{j=\ell+1}^k \tilde{u}^j$ with $\tilde{u}^j \in V(\lambda_j)$ such that $\sum_{j=\ell+1}^k C_{\tilde{\theta}_{\ell+1}}^j e^{J_j t} \tilde{u}^j \geq h_2$, $\forall t \geq 0$. Notice

that $\lambda_\ell > \lambda_j$ for all $j = \ell + 1, \dots, k$. Hence, $e^{-\lambda_\ell t} \|\sum_{j=\ell}^k C_{\theta_\ell}^j e^{J_j t} \tilde{u}^j\|_2$ is bounded on $[0, \infty)$ as \tilde{u}^j 's are fixed, namely, there exists a positive scalar ρ satisfying $e^{-\lambda_\ell t} \|\sum_{j=\ell}^k C_{\theta_\ell}^j e^{J_j t} \tilde{u}^j\|_2 \leq \rho$, $\forall t \geq 0$.

Moreover, by Proposition 8 and $P(C_{\theta_\ell}^\ell, v^\ell) \succ 0$, we deduce the existence of $\tilde{u}^\ell \in V(\lambda_\ell)$ such that $C_{\theta_\ell}^\ell e^{J_\ell t} \tilde{u}^\ell \geq e^{\lambda_\ell t} (h_1 + \rho \mathbf{1})$, where $\mathbf{1} = (1, \dots, 1)^T$. Consequently, $C_{\theta_\ell}^\ell e^{J_\ell t} \tilde{u}^\ell + \sum_{j=\ell+1}^k C_{\theta_\ell}^j e^{J_j t} \tilde{u}^j \geq e^{\lambda_\ell t} h_1 \geq h_1$, $\forall t \geq 0$. Furthermore, it is noted, by the definition of θ_j and the echelon structure of \tilde{P} , that $P(C_{\theta_j}^\ell, v^\ell) \equiv 0$ for all $j = \ell + 1, \dots, k$. Hence, we obtain, via Proposition 8, that $C_{\theta_j}^\ell e^{J_\ell t} \tilde{u}^\ell \equiv 0$ for all $j = \ell + 1, \dots, k$. As a result, for any $i = \ell + 1, \dots, k$, we have $C_{\theta_i}^\ell e^{J_\ell t} \tilde{u}^\ell + \sum_{j=\ell+1}^k C_{\theta_i}^j e^{J_j t} \tilde{u}^j =$

$\sum_{j=\ell+1}^k C_{\theta_i}^j e^{J_j t} \tilde{u}^j$, $\forall t \geq 0$. Letting $u = \bigoplus_{j=\ell}^k \tilde{u}_i \in \bigoplus_{j=\ell}^k V(\lambda_j)$, it is easy to see from the above development that

$$\sum_{j=\ell}^k C_{\theta_i}^j e^{J_j t} \tilde{u}^j = \begin{pmatrix} C_{\theta_\ell} e^{J_\ell t} \tilde{u}^\ell \\ \text{-----} \\ C_{\tilde{\theta}_{\ell+1}} e^{J_{\ell+1} t} \tilde{u}^{\ell+1} \end{pmatrix} = \begin{pmatrix} \sum_{j=\ell}^k C_{\theta_\ell}^j e^{J_j t} \tilde{u}^j \\ \text{-----} \\ \sum_{j=\ell+1}^k C_{\theta_{\ell+1}}^j e^{J_j t} \tilde{u}^j \\ \vdots \\ \sum_{j=\ell}^k C_{\theta_k}^j e^{J_j t} \tilde{u}^j \end{pmatrix} = \begin{pmatrix} \sum_{j=\ell}^k C_{\theta_\ell}^j e^{J_j t} \tilde{u}^j \\ \text{-----} \\ \sum_{j=\ell+1}^k C_{\theta_{\ell+1}}^j e^{J_j t} \tilde{u}^j \\ \vdots \\ \sum_{j=\ell+1}^k C_{\theta_k}^j e^{J_j t} \tilde{u}^j \end{pmatrix} \geq \begin{pmatrix} h_1 \\ \text{---} \\ h_2 \end{pmatrix} = h$$

for all $t \geq 0$. This shows that u is the desired vector and thus completes the proof of the claim.

Letting $i = 1$, then for $h \in \mathbb{R}_{++}^{|\theta|}$ with $h_\alpha = b^+ > 0$, we obtain $u \equiv \bigoplus_{j=1}^k u^j$ with $u^j \in V(\lambda_j)$ satisfying

the condition stated in the claim. Note that $C_{\bar{\theta}}^j e^{J_j t} u^j \equiv 0, \forall t$ for all $j = 1, \dots, k$, where $\bar{\theta} \subseteq \beta$.

Therefore, appropriately expanding u to $u^* = (u^T, 0)^T \in \mathbb{R}^n$ by adding the zero subvectors, we

have $C_\alpha e^{At} u^* = \sum_{j=1}^k C_\alpha^j e^{J_j t} u^j \geq b^+$ and $C_\beta e^{At} u^* = \sum_{j=1}^k C_\beta^j e^{J_j t} u^j \geq 0$ for all $t \geq 0$. This shows

$u^* \in \mathcal{A}$.

(S2) $\alpha \neq \emptyset, \gamma \neq \emptyset$. If, in addition to $P^+ \succ 0$ and $P^0 \succcurlyeq 0$, $P^- \succcurlyeq 0$ holds true, then we

have $\tilde{P} \succ 0$ and $(P(C_\theta^1, v^1), \dots, P(C_\theta^k, v^k)) \equiv 0$, where \tilde{P} is defined as before with $\theta \subseteq \beta \cup \gamma$.

Hence we can apply the similar argument as in (S1) to obtain $u^* \in \mathbb{R}^n$ such that $C_\alpha e^{At} u^* \geq b^+$,

$C_\beta e^{At} u^* \geq 0$, and $C_\gamma e^{At} u^* \geq 0 \geq b^-$ for all $t \geq 0$. Thus $u^* \in \mathcal{A}$. Now suppose $P^- \not\succeq 0$ but the

condition (7) holds. Let $\bar{\psi} \equiv \{1, \dots, m\} \setminus \psi$, where ψ is defined in (7). Therefore $\bar{\psi} = \alpha \cup \beta \cup (\gamma \setminus \psi)$.

Hence, $P_{\bar{\psi}} \equiv (P(C_{\bar{\psi}}^1, v^1), \dots, P(C_{\bar{\psi}}^k, v^k)) \succcurlyeq 0$. By (S1) shown above, there exists $u = \bigoplus_{j=1}^k u^j$ with

$u^j \in V(\lambda_j)$ such that

$$\sum_{j=1}^k C_\alpha^j e^{J_j t} u^j \geq b^+, \quad \sum_{j=1}^k C_\beta^j e^{J_j t} u^j \geq 0, \quad \sum_{j=1}^k C_{\gamma \setminus \psi}^j e^{J_j t} u^j \geq 0 \geq b_{\gamma \setminus \psi}, \quad \forall t \geq 0$$

Particularly, since $\sum_j C_\phi^{kj} \mathcal{L}^0(v^{kj}) = \sum_j C_\phi^{kj} v^{kj} \geq b_\phi$, we deduce, via the proof of Proposition 8

and Remark 9, that the positive coefficient μ_0 in u^k can be taken as one. Therefore, using the

lexicographical relation and $\lambda_k = 0$, we have

$$\sum_{i=1}^k C_\psi^i e^{J_i t} u^i = C_\psi^k e^{J_k t} u^k = e^{\lambda_k t} \mu_0 \sum_j C_\psi^{kj} \mathcal{L}^0(v^{kj}) = \sum_j C_\psi^{kj} v^{kj} \geq b_\psi, \quad \forall t \geq 0.$$

Consequently, appropriately expanding u to a vector $u^* \in \mathbb{R}^n$ by adding the zero subvectors, we

obtain $u^* \in \mathcal{A}$.

“Necessity”. Let $x^* \in \mathcal{A}$, i.e., $Ce^{At} x^* \geq b$ for all $t \geq 0$. Since $\alpha \neq \emptyset$, $Ce^{At} x^*$ is not identically zero, i.e., $x^* \notin \overline{O}(C, A)$, where $\overline{O}(C, A)$ denotes the unobservable subspace of the pair (C, A) . We deduce from (5) that

$$Ce^{At} x^* = \sum_i \left\{ e^{\lambda_i t} \sum_{k=0}^{\hat{n}_i-1} \frac{t^k}{k!} \left(\sum_j C^{ij} \mathcal{L}^k(v^{ij}) + \sum_s u_s^{ik} \sin(\omega_s^{ik} t + \theta_s^{ik}) \right) \right\}, \quad (11)$$

where λ_i 's are the real parts of the eigenvalues of A , \hat{n}_i is the largest order of the Jordan blocks

associated with (possibly complex) eigenvalues with the real part λ_i , and (possibly zero) $u_s^{ik} \in \mathbb{R}^m$,

$\omega_s^{ik} \in \mathbb{R}_{++}$, $\theta_s^{ik} \in \mathbb{R}$ depends on C, A, x^* only. It should be noted that if λ_i is such that the

matrix A has no Jordan block associated with a real eigenvalue λ_i (namely, λ_i corresponds to a

strictly complex eigenvalue of A), then $\sum_j C^{ij} \mathcal{L}^k(v^{ij}) = 0$ for all k in (11). Let the real parts of the

eigenvalues of A be labeled in a descending order, i.e., $\lambda_1 > \lambda_2 > \dots > \lambda_k \geq 0 > \lambda_{k+1} > \dots > \lambda_q$.

For each nonnegative real part λ_i , define the vector tuple $P(C^i, v^i)$ as follows: if λ_i corresponds a

strictly complex eigenvalue of A , then we set $P(C^i, v^i) \equiv 0$; otherwise, i.e., λ_i is a real eigenvalue

of A ,

$$P(C^i, v^i) \equiv \left(\sum_j C^{ij} \mathcal{L}^{\hat{n}_i-1}(v^{ij}), \sum_j C^{ij} \mathcal{L}^{\hat{n}_i-2}(v^{ij}), \dots, \sum_j C^{ij} \mathcal{L}^0(v^{ij}) \right), \quad i = 1, \dots, k$$

where $v^i = \bigoplus_j v^{ij} \in V(\lambda_i)$ is the subvector in x^* corresponding to λ_i . We show below that the tuples $P^+ \equiv (P(C_\alpha^1, v^1), \dots, P(C_\alpha^k, v^k))$, $P^0 \equiv (P(C_\beta^1, v^1), \dots, P(C_\beta^k, v^k))$, and $P^- \equiv (P(C_\gamma^1, v^1), \dots, P(C_\gamma^k, v^k))$ satisfy the conditions (a) and (b) in the theorem. For notational simplicity, let P_i^+, P_i^0, P_i^- denote the i th row of P^+, P^0 , and P^- respectively.

(N1) To prove $P^+ \succ 0$, it is sufficient to show that for each $\ell \in \alpha$, $P_\ell^+ \succ 0$, i.e., $P_\ell^+ \neq 0$ and the leading entry in P_ℓ^+ is positive. Let the real pair (λ, p) represent the ‘‘mode’’ of $e^{\lambda t} t^p$ in $C_\ell e^{At} x^*$, where p is a nonnegative integer, and let $(\tilde{\lambda}, \tilde{p})$ represent the largest non-vanishing mode in $C_\ell e^{At} x^*$ in the sense that $(\tilde{\lambda}, \tilde{p}) \succ (\lambda_i, p_i)$ for any other non-vanishing mode defined by (λ_i, p_i) . Since $C_\ell e^{At} x^* \geq b_\ell > 0, \forall t \geq 0$, we must have $(\tilde{\lambda}, \tilde{p}) \succcurlyeq (0, 0)$ because otherwise, $C_\ell e^{At} x^* \rightarrow 0$ as $t \rightarrow \infty$, a contradiction. Since $\tilde{\lambda} \geq 0, \tilde{\lambda} = \lambda_r$ for some $r \in \{1, \dots, k\}$. Furthermore, we

deduce from (11) that $C_\ell e^{At} x^*$ tends to $e^{\lambda_r t} \frac{t^{\tilde{p}}}{\tilde{p}!} \left[\rho_0 + \sum_{s=1}^{k_\ell} \rho_s \sin(\omega_s t + \theta_s) \right]$ as $t \rightarrow +\infty$, where

$(\rho_0, \rho_1, \dots, \rho_{k_\ell}) \neq 0$ depends on C_ℓ, A, x^* only and $\rho_0 = \sum_j C_\ell^{rj} \mathcal{L}^{\tilde{p}}(v^{rj})$. Since (λ_r, \tilde{p}) denotes the largest non-vanishing mode in $C_\ell e^{At} x^*$, all the terms before ρ_0 in P_ℓ^+ must be zero. Therefore, ρ_0 is the first nonzero entry in P_ℓ^+ . Since $b_\ell > 0$, we deduce $\rho_0 > 0$ in view of (a) of Proposition 11 as desired.

(N2) To prove $P^0 \succcurlyeq 0$, we show that $P_{i_\ell}^0 \succcurlyeq 0$ for each $\ell \in \beta$, where i_ℓ is the row index of P^0 corresponding to $\ell \in \beta$. Suppose, in contrast, that there exists $\ell \in \beta$ such that the first nonzero element in $P_{i_\ell}^0$ is negative which corresponds to the mode in $C_\ell e^{At} x^*$ defined by (λ_r, \tilde{p}) . Hence, $(\lambda_r, \tilde{p}) \succcurlyeq (0, 0), \sum_j C_\ell^{rj} \mathcal{L}^{\tilde{p}}(v^{rj}) < 0, \sum_j C_\ell^{rj} \mathcal{L}^p(v^{rj}) = 0$ for $\tilde{p} < p \leq \bar{n}_r - 1$, and $P(C_\ell^i, v^i) = 0$ for all $i = 1, \dots, r + 1$. Let $(\hat{\lambda}, \hat{p})$ represent the largest non-vanishing mode in $C_\ell e^{At} x^*$. Therefore, $(\hat{\lambda}, \hat{p})$ exists and satisfies $(\hat{\lambda}, \hat{p}) \succcurlyeq (\lambda_r, \tilde{p}) \succcurlyeq (0, 0)$. Moreover, $C_\ell e^{At} x^*$ tends

to $e^{\hat{\lambda} t} \frac{t^{\hat{p}}}{\hat{p}!} \left[\rho_0 + \sum_{s=1}^{k_\ell} \rho_s \sin(\omega_s t + \theta_s) \right]$ as $t \rightarrow +\infty$, where $(\rho_0, \rho_1, \dots, \rho_{k_\ell}) \neq 0$. Notice that either

$\rho_0 = 0$ when $(\hat{\lambda}, \hat{p}) \succ (\lambda_r, \tilde{p})$ or $\rho_0 < 0$ when $(\hat{\lambda}, \hat{p}) = (\lambda_r, \tilde{p})$. However, this contradicts (a) of Proposition 11 since $b_\ell = 0$.

(N3) To show the implication (7) holds for P^- , define the (nonempty) index set $\vartheta \equiv \{j \in \gamma \mid \text{the } j\text{th row of } P^- \text{ has a negative leading entry}\}$. For each $\ell \in \vartheta$, let the pair $(\lambda_r, \tilde{p}) \succcurlyeq (0, 0)$ correspond to the mode in $C_\ell e^{At} x^*$ associated with the negative leading entry in $P_{i_\ell}^-$ (again i_ℓ corresponds to the $\ell \in \vartheta$), and let $(\hat{\lambda}, \hat{p}) \succcurlyeq (\lambda_r, \tilde{p})$ correspond to the largest non-vanishing mode in $C_\ell e^{At} x^*$. We claim that $(\hat{\lambda}, \hat{p}) = (\lambda_r, \tilde{p})$. Suppose not. Then $(\hat{\lambda}, \hat{p}) \succ (0, 0)$ and $C_\ell e^{At} x^*$ tends to $e^{\hat{\lambda} t} \frac{t^{\hat{p}}}{\hat{p}!} \left[\rho_0 + \sum_{s=1}^{k_\ell} \rho_s \sin(\omega_s t + \theta_s) \right]$ as $t \rightarrow +\infty$ with $\rho_0 = 0$ and $(\rho_1, \dots, \rho_{k_\ell}) \neq 0$, a contradiction to (a)

of Proposition 11. Hence $(\hat{\lambda}, \hat{p}) = (\lambda_r, \tilde{p})$ and $C_\ell e^{At} x^*$ is tends to $e^{\lambda_r t} \frac{t^{\tilde{p}}}{\tilde{p}!} \left[\tilde{\rho}_0 + \sum_{s=1}^{k_\ell} \tilde{\rho}_s \sin(\omega_s t + \theta_s) \right]$

for $t \rightarrow +\infty$ with $\tilde{\rho}_0 < 0$. We thus deduce $(\lambda_r, \tilde{p}) = (0, 0)$ and $\lambda_r = \lambda_k = 0$ by (b) of Proposition 11, namely, the negative leading entry in $P_{i_\ell}^-$ appears in the last column of P^- . This shows that ϑ equals to ψ defined in (7) and thus $(P(C_{\gamma \setminus \psi}^1, v^1), \dots, P(C_{\gamma \setminus \psi}^k, v^k)) \succcurlyeq 0$. Moreover, we obtain $\tilde{\rho}_0 \equiv \sum_j C_\ell^{kj} v^{kj} \geq b_\ell$ for each $\ell \in \psi$ from (c) of Proposition 11. Therefore, $\sum_j C_\psi^{kj} v^{kj} \geq b_\psi$ holds. Following the similar argument, we also have, for each $\ell \in \phi \subseteq \alpha, C_\ell e^{At} x^*$

tends to $\sum_j C_\ell^{kj} v^{kj} + \sum_{s=1}^{k_\ell} \tilde{\rho}_s \sin(\omega_s t + \theta_s)$ as $t \rightarrow +\infty$. This yields, via (c) of Proposition 11, that $\sum_j C_\ell^{kj} v^{kj} \geq b_\ell > 0$. As a result, $\sum_j C_\phi^{kj} v^{kj} \geq b_\phi$.

Finally we remove each nonnegative λ_i corresponding to a strictly complex eigenvalue of A from

$(\lambda_1, \dots, \lambda_k)$ and its associated zero block from P^+ , P^0 , P^- . Since the index set α is nonempty, there exists $\ell \in \alpha$ such that the long-time dominating mode of $C_\ell e^{At} x^*$ has a positive ρ_0 as shown in (N1). This implies that at least one of $\lambda_i \geq 0, i = 1, \dots, k$ is a real eigenvalue of A . For each remaining real eigenvalue $\lambda_i \geq 0$, let \bar{n}_i be the largest order of the Jordan blocks associated with the real λ_i . If the integer \hat{n}_i , the largest order of the Jordan blocks associated with (possibly complex) eigenvalues with the real part λ_i , is greater than \bar{n}_i , then we can replace $P(C^i, v^i)$ by

$$P(C^i, v^i) \equiv \left(\sum_j C^{ij} \mathcal{L}^{\bar{n}_i-1}(v^{ij}), \sum_j C^{ij} \mathcal{L}^{\bar{n}_i-2}(v^{ij}), \dots, \sum_j C^{ij} \mathcal{L}^0(v^{ij}) \right),$$

since $\sum_j C^{ij} \mathcal{L}^p(v^{ij}) = 0$ for $p \geq \bar{n}_i$. We thus obtain the nonnegative real eigenvalues $\tilde{\lambda}_1 > \dots > \tilde{\lambda}_k \geq 0$ and the tuples $\tilde{P}^+, \tilde{P}^0, \tilde{P}^-$. It is easy to verify that removing the zero blocks from P^+, P^0, P^- does not change the lexicographical relations or the implications proved in (N1–N3) for $\tilde{P}^+, \tilde{P}^0, \tilde{P}^-$, i.e., $\tilde{P}^+ \succ 0, \tilde{P}^- \succcurlyeq 0$ and the implication (7) holds for \tilde{P}^- . This shows that nonnegative eigenvalues $\tilde{\lambda}_i$ and $v^i \in V(\tilde{\lambda}_i)$ satisfy the conditions (a)–(b). \square

7.2 Generalized Sublinear Functions

Proof of Lemma 22. Since f is nontrivial, there exists a nonzero $x^* \in \text{dom} f$ such that $f(x^*) \neq 0$. Hence, we have $f(0) = f(0 \cdot x^*) = p(0)f(x^*)$, which yields $f(0) < +\infty$. Furthermore, noting that $f(x^*) = f(\lambda \cdot \frac{1}{\lambda} \cdot x^*)$ for all $\lambda > 0$, we have $f(x^*) = p(\lambda)p(\frac{1}{\lambda})f(x^*)$. This shows that $p(\lambda)p(\frac{1}{\lambda}) = 1, \forall \lambda > 0$ and $p^2(1) = 1$ by letting $\lambda = 1$. Since $p(\lambda) \geq 0, \forall \lambda \geq 0$, we have $p(1) = 1$. The increasing property of p further implies that there exists $\lambda^* > 1$ with $p(\lambda^*) \neq 1$. Observing $f(0) = p(\lambda)f(0), \forall \lambda \geq 0$, we obtain $f(0) = p(\lambda^*)f(0)$ and thus $f(0) = 0$ by using $f(0) < +\infty$ shown before. Now let $\lambda = 0$. For any $x \in \text{dom} f$ with $f(x) \neq 0$, we have $f(0) = p(0)f(x)$. Therefore, $p(0) = 0$. This yields $p(\lambda) > 0, \forall \lambda > 0$ and thus $p(\frac{1}{\lambda}) = \frac{1}{p(\lambda)}$ for all $\lambda > 0$. \square

7.3 Switching Properties of PASs

Associated with the PAS (9), we define a reverse-time system as follows: for any given terminal time $T > 0$, let $x^r(t) \equiv x(T - t)$ and $x^r(0) = x(T)$. Then we have

$$\dot{x}^r = -A_i x^r - d_i, \quad \text{if } x^r \in \mathcal{P}_i. \quad (12)$$

It is easy to show that the reverse-time system (12) remains a PAS. We call the PAS (9) *forward-time non-Zeno* if for any $x^0 \in \mathbb{R}^n$, there exist a scalar $\varepsilon > 0$ and a convex polyhedron \mathcal{P}_i such that $x(t, x^0) \in \mathcal{P}_i$ for all $t \in [0, \varepsilon]$. Similarly, the PAS (9) is *backward-time non-Zeno* if the associated reverse-time PAS is forward-time non-Zeno. If a PAS is both forward-time and backward-time non-Zeno, then we call it *non-Zeno*.

To characterize the local solution properties of the PAS, we introduce the vector $z \in \mathbb{R}^{n+n+m}$ and write the i th mode of the PAS in the following equivalent homogeneous form

$$\dot{z} = \tilde{A}_i z, \quad \text{if } z \in \tilde{\mathcal{P}}_i \equiv \{z \mid \tilde{C}_i z \geq 0\},$$

where

$$\tilde{A}_i = \begin{bmatrix} A_i & I & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \in \mathbb{R}^{(2n+m) \times (2n+m)}, \quad \tilde{C}_i = [C_i \quad 0 \quad I] \in \mathbb{R}^{m \times (2n+m)}, \quad z = \begin{bmatrix} x \\ d_i \\ -b_i \end{bmatrix} \in \mathbb{R}^{n+n+m}.$$

For the i th mode of the PAS, let

$$\mathcal{Y}_i \equiv \left\{ x \in \mathbb{R}^n \mid (\tilde{C}_i z, \tilde{C}_i \tilde{A}_i z, \dots, \tilde{C}_i \tilde{A}_i^{2n+m-1} z) \succcurlyeq 0, \quad z = \begin{bmatrix} x \\ d_i \\ -b_i \end{bmatrix} \right\}.$$

For any $x^0, x(t, x^0) \in \mathcal{P}_i$ for all $t \geq 0$ sufficiently small if and only if $x^0 \in \mathcal{Y}_i$. Given $x^0 \in \mathbb{R}^n$, define two index sets $\mathcal{I}(x^0) \equiv \{i \mid x^0 \in \mathcal{P}_i\}$ and $\mathcal{J}(x^0) \equiv \{i \mid x^0 \in \mathcal{Y}_i\}$. It is clear that $\mathcal{J}(x^0) \subseteq \mathcal{I}(x^0)$. Similarly, we can define $\mathcal{J}^r(x^0)$ for the associated reverse-time system. We call a PAS of *simple switching property* if for any state trajectory $x(t, x^0)$, $\mathcal{J}(x(t_*, x^0)) = \mathcal{J}^r(x(t_*, x^0))$ at any non-switching time t_* whose definition is given in Section 4.1.

The non-Zeno and simple switching properties of the PASs can be proved via the similar techniques as those in [11, 34] for the CLSs. To be self-contained, we provide concise proofs with an emphasis on different arguments due to the presence of affine dynamics and polyhedrons for the PASs (in contrast to linear dynamics and polyhedral cones for the CLSs) in this section. The interested reader may refer to [11, 34] for more technical details.

Proof of non-Zenoness. It is sufficient to show the forward-time non-Zenoness of the PAS, which is equivalent to the condition that for any $x^0, x^0 \in \mathcal{Y}_i$ for some i . For any $x^0 \in \mathbb{R}^n$ and $i \in \mathcal{I}(x^0)$, let $z^{i0} \equiv ((x^0)^T, d_i^T, -b_i^T)^T$. We collect the following facts that can be proved in the similar way as in [11]:

- (a) $x^0 \in \mathcal{Y}_i$ if and only if $\sum_{k=0}^{2n+m-1} t^k \tilde{A}_i^k z^{i0} \in \tilde{\mathcal{P}}_i$ for all $t \geq 0$ sufficiently small;
- (b) $\tilde{A}_i z^{i0} = \tilde{A}_j z^{j0}$ for any $i, j \in \mathcal{I}(x^0)$;
- (c) $\tilde{A}_i^k z^{i0} = \tilde{A}_j^k z^{j0}$ for all $k \geq 0$ and any $i, j \in \mathcal{J}(x^0)$;
- (d) for a vector-valued polynomial $p(t)$ with $p(0) = x^0$, there exists $i \in \mathcal{I}(x^0)$ such that $p(t) \in \mathcal{P}_i$ for all $t \geq 0$ sufficiently small.

For a given x^0 , let $\mathcal{I}_0 \equiv \mathcal{I}(x^0)$. Define $z^0 \equiv z^{i0}$ and $z^1 \equiv \tilde{A}_i z^0$ for any $i \in \mathcal{I}_0$. We deduce from (b) that z^0 and z^1 are unique. Now let $\mathcal{I}_1 \equiv \{i \in \mathcal{I}_0 \mid z^0 + t z^1 \in \tilde{\mathcal{P}}_i \text{ for } t > 0 \text{ sufficiently small}\}$. We claim: (i) \mathcal{I}_1 is nonempty; (ii) $\mathcal{I}_1 \subseteq \mathcal{I}_0$; and (iii) $\tilde{A}_i z^1 = \tilde{A}_j z^1$ for $i, j \in \mathcal{I}_1$. Indeed, since $\mathcal{I}_1 = \{i \in \mathcal{I}_0 \mid (C x^0 - b_i) + t C_i (A_i x^0 + d_i) \geq 0 \text{ for } t > 0 \text{ sufficiently small}\} = \{i \in \mathcal{I}_0 \mid x^0 + t(A_i x^0 + d_i) \in \mathcal{P}_i \text{ for } t > 0 \text{ sufficiently small}\}$, (i) follows from (d). Moreover, (ii) is obvious. To see (iii), it is clear from the above analysis that $\mathcal{I}_1 = \mathcal{I}([x^0 + t(A_i x^0 + d_i)])$ for $t \geq 0$ sufficiently small. Hence, for $i, j \in \mathcal{I}_1$, we have, via (b), that $\tilde{A}_i(z^0 + t z^1) = \tilde{A}_j(z^0 + t z^1)$ for all $t \geq 0$ small. This gives rise to (iii) and shows the claim. Inductively, we define $\mathcal{I}_\ell \equiv \{i \in \mathcal{I}_0 \mid z^0 + t z^1 + \dots + t^\ell z^\ell \in \tilde{\mathcal{P}}_i \text{ for } t > 0 \text{ sufficiently small}\}$, where $z^j = \tilde{A}_i z^{j-1}$ for some $i \in \mathcal{I}_{j-1}$, $j = 1, \dots, \ell$. We claim for each $\ell \geq 1$: (i) \mathcal{I}_ℓ is nonempty; (ii) $\mathcal{I}_\ell \subseteq \mathcal{I}_{\ell-1}$; and (iii) $\tilde{A}_i z^\ell = \tilde{A}_j z^\ell$ for $i, j \in \mathcal{I}_\ell$. The case of $\ell = 1$ has been proven. Now suppose the claims holds for all $\ell = 1, \dots, p$, where $p \geq 1$, and consider $\ell = p + 1$. The nonemptiness of \mathcal{I}_{p+1} is due to the similar reason as for $\ell = 1$. To see $\mathcal{I}_\ell \subseteq \mathcal{I}_{\ell-1}$, it is observed that for any $i \in \mathcal{I}_\ell$, $\tilde{C}_i(z^0 + t z^1 + \dots + t^\ell z^\ell) \geq 0$ for $t \geq 0$ sufficiently small. This shows that $(\tilde{C}_i z^0, \tilde{C}_i z^1, \dots, \tilde{C}_i z^\ell) \succcurlyeq 0$, which further implies $(\tilde{C}_i z^0, \tilde{C}_i z^1, \dots, \tilde{C}_i z^{\ell-1}) \succcurlyeq 0$. The latter yields $\tilde{C}_i(z^0 + t z^1 + \dots + t^{\ell-1} z^{\ell-1}) \geq 0$ for $t \geq 0$ sufficiently small, or equivalently $i \in \mathcal{I}_{\ell-1}$. Finally, (iii) holds because the similar argument for $\ell = 1$. To complete the non-Zeno proof, letting $\ell = 2n + m - 1$ and using (a), we reach the desired result. \square

Proof of the simple switching property. We prove that a time $t_* \geq 0$ is a non-switching time along a state trajectory $x(t, x^0)$ if and only if $\mathcal{J}(x(t_*, x^0)) = \mathcal{J}^r(x(t_*, x^0))$. Let $x^* \equiv x(t_*, x^0)$. It can be shown using the same argument in [11, Proposition 3.11] that t_* is a non-switching time if and only if $\mathcal{J}(x^*) \cap \mathcal{J}^r(x^*)$ is nonempty. Therefore, the “if” part is a direct consequence of the fact that $\mathcal{J}(x^*)$ is nonempty due to the non-Zenoness proved previously. For the “only if” part, since t_* is a non-switching time, $\mathcal{J}(x^*) \cap \mathcal{J}^r(x^*)$ is nonempty. Let $j \in \mathcal{J}(x^*) \cap \mathcal{J}^r(x^*)$. It is sufficient to show $\mathcal{J}^r(x^*) \subseteq \mathcal{J}(x^*)$ since $\mathcal{J}^r(x^*) \supseteq \mathcal{J}(x^*)$ can be proved in the similar way via the reverse-time system. Suppose this is not the case. Then there exists $i \in \mathcal{J}^r(x^*)$ but $i \notin \mathcal{J}(x^*)$.

Hence $x(t, x^0) \in \mathcal{P}_i \cap \mathcal{P}_j$ on $[t_* - \varepsilon, t_*]$ for some $\varepsilon > 0$, namely, $x(t, x^0) \in \mathcal{P}_i \cap \{x \mid (C_i x - b_i)_\alpha = 0\} = \mathcal{P}_j \cap \{x \mid (C_j x - b_j)_\beta = 0\}$ on $[t_* - \varepsilon, t_*]$ for nonempty index sets α and β . This implies that $((x^*)^T, d_j^T, -b_j^T) \in \overline{O}((\tilde{C}_j)_\beta, \tilde{A}_j)$, which further shows that $(C_j x(t, x^0) - b_j)_\beta = 0$ on $[t_*, t_* + \varepsilon]$. Hence, $x(t, x^0) \in \mathcal{P}_j \cap \{x \mid (C_j x - b_j)_\beta = 0\} = \mathcal{P}_i \cap \mathcal{P}_j$ on $[t_*, t_* + \varepsilon]$. Consequently, $x(t, x^0) \in \mathcal{P}_i$ on $[t_*, t_* + \varepsilon]$ and thus $i \in \mathcal{J}(x^*)$. This is a contradiction. \square

Proof of Lemma 25. The proof is similar to that of [34, Proposition 10]. To see the “if” part, it is observed that $x(t, x^0) \notin \bigcup_i \mathcal{E}_i$ at each $t \geq 0$ for any non-equilibrium x^0 . Since $\mathcal{A}_i = \mathcal{E}_i$ for all i , we have $x(t, x^0) \notin \bigcup_i \mathcal{A}_i$ for all $t \geq 0$. This shows that there are infinite mode switchings along $x(t, x^0)$ in the positive time direction in view of the simple switching property. On the other hand, suppose the PAS has infinite mode switchings but $\mathcal{A}_i \neq \mathcal{E}_i$ for some i . Since $\mathcal{E}_i \subseteq \mathcal{A}_i$, there exists $x^0 \in \mathcal{A}_i \subseteq \mathcal{X}_i$ but $x^0 \notin \mathcal{E}_i$. This implies that $x(t, x^0)$ has no switching for all $t \geq 0$, via the simple switching property. This is a contradiction as x^0 is not an equilibrium. \square

Acknowledgments. The author thanks the associate editor and the three reviewers for their constructive comments that have significantly improved the presentation of the paper.

References

- [1] A. ABATE, A. TIWARI, AND S. SASTRY. Box invariance for biologically-inspired dynamical systems. *Proceedings of the 46th IEEE Conference on Decision and Control*, New Orleans, LA, 2007.
- [2] E.J. ANDERSON AND P. NASH. *Linear Programming in Infinite-Dimensional Spaces*. Wiley, New York, 1987.
- [3] V.I. ARNOLD. *Mathematical Methods of Classical Mechanics*. 2nd Edition, Springer, New York, 1997.
- [4] J.P. AUBIN. *Viability Theory*. Birkhäuser, Boston, 1991.
- [5] G. BATT, C. BELTA, AND R. WEISS. Temporal logic analysis of gene networks under parameter uncertainty. *IEEE Trans. on Automatic Control*, Vol.53(1), Special issue on systems biology, pp.215–228, 2008.
- [6] C. BELTA AND L.C.G.J.M. HABETS. Controlling of a class of nonlinear systems on rectangles. *IEEE Trans. on Automatic Control*, Vol.51(11), pp.1749–1759, 2006.
- [7] A. BERMAN, M. NEUMANN, AND R.J. STERN. *Nonnegative Matrices in Dynamical Systems*. John Wiley & Sons, New York, 1989.
- [8] V. BLONDEL, AND J.N. TSITSIKLIS. A survey of computational complexity results in systems and control. *Automatica*, Vol.36, pp.1249–1274, 2000.
- [9] J. BOCHNAK, M. COSTE, AND M.-F. ROY. *Real Algebraic Geometry*. Springer, 1998.
- [10] J.M. BORWEIN AND A.S. LEWIS. *Convex Analysis and Nonlinear Optimization*. 2nd Edition, Springer, 2006.
- [11] M.K. ÇAMLIBEL, J.S. PANG, AND J. SHEN. Conewise linear systems: non-Zenoness and observability. *SIAM Journal on Control and Optimization*, Vol.45(6), pp.1769–1800, 2006.
- [12] E.B. CASTELAN AND J.C. HENNET. On invariant polyhedra of continuous-time linear systems. *IEEE Trans. on Automatic Control*, Vol.38(1), pp.1680–1685, 1993.
- [13] P. CHENG, G.J. PAPPAS, AND V. KUMAR. Decidability of motion planning with differential constraints. *Proceedings of 2007 IEEE International Conference on Robotics and Automation*, pp.1826–1831, Roma, Italy, 2007.

- [14] H. DE JONG. Modeling and simulation of genetic regulatory systems: a literature review. *Journal of Computational Biology*, Vol.9(1), pp.67–103, 2002.
- [15] F. FACCHINEI AND J.S. PANG. *Finite-Dimensional Variational Inequalities and Complementarity Problems*. Springer-Verlag, New York, 2003.
- [16] KY FAN. On systems of linear inequalities. *Linear Inequalities and Related Systems*, Annals of Mathematical Studies, No.38, Edited by H.W. Kuhn and A.W. Tucker, pp.99-156, Princeton University Press, 1956.
- [17] R. GHOSH AND C.J. TOMLIN. Symbolic reachable set computation of piecewise affine hybrid automata and its application to biological modelling: Delta-notch protein signalling. *IEE Systems Biology*, Vol.1(1), pp.170–183, June, 2004.
- [18] L.C.G.J.M. HABETS, P.J. COLLINS, AND J.H. VAN SCHUPPEN. Reachability and control synthesis for piecewise-affine hybrid systems on simplices. *IEEE Trans. on Automatic Control*, Vol.51(6), pp.938–947, 2006.
- [19] J.-B. HIRIART-URRUTY AND C. LEMARÉCHAL. *Foundamentals of Convex Analysis*. Grundlehren Text Editions, Springer-Verlag, Berlin, 2001.
- [20] H. KHALIL. *Nonlinear Systems*. 2nd Edition, Prentice Hall, 1996.
- [21] G. LAFFERRIERE, G.J. PAPPAS, AND S. SASTRY. O-minimal hybrid systems. *Mathematics of Control, Signals, and Systems*, Vol.13, pp.1–21, 2000.
- [22] G. LAFFERRIERE, G.J. PAPPAS, AND S. YOVINE. Symbolic reachability computation for families of linear vector fields. *Journal of Symbolic Computation*, Vol.32, pp.231–253, 2001.
- [23] G. LAFFERRIERE, G.J. PAPPAS, AND S. YOVINE. A new class of decidable hybrid systems. *Hybrid Systems: Computation and Control*, LNCS 1569, pp.137–151, Berlin, Springer, 1999.
- [24] M. PACHTER AND D.H. JACOBSON. Observability with a conic observation set. *IEEE Transactions on Automatic Control*, Vol.24(4), pp.632–633, 1979.
- [25] J.S. PANG AND J. SHEN. Strongly regular differential variational systems. *IEEE Transactions on Automatic Control*, Vol.52(2), pp.242–255, 2007.
- [26] J.S. PANG AND D.E. STEWART. Differential variational inequalities. *Mathematical Programming, Series A*, Vol.113(2), pp.345–424, 2008.
- [27] P.A. PARRILO. Semidefinite programming relaxations for semialgebraic problems. *Mathematical Programming Series B*, Vol.96(2), pp.293–320, 2003.
- [28] S. PRAJNA, A. JADBABAIE, AND G.J. PAPPAS. A framework for worst-case and stochastic safety verification using barrier certificates. *IEEE Trans. on Automatic Control*, Vol.52(8), pp.1415–1428, 2007.
- [29] M. PRANDINI AND J. HU. A stochastic approximation method for reachability computations. *Stochastic Hybrid Systems: Theory and Safety Applications*, Lecture Notes in Control and Information Sciences, Springer-Verlag, 2006.
- [30] S.H. SAPERSTONE AND J.A. YORKE. Controllability of linear oscillatory systems using positive controls. *SIAM Journal on Control*, Vol.9(2), pp.253–262, 1971.
- [31] S. SCHOLTES. *Introduction to Piecewise Differentiable Equations*. Habilitation thesis, Institut für Statistik und Mathematische Wirtschaftstheorie, Universität Karlsruhe, 1994.
- [32] J.M. SCHUMACHER. Complementarity systems in optimization. *Mathematical Programming, Series B*, Vol. 101, pp. 263–295, 2004.

- [33] J. SHEN. Positive invariance and observability of conewise linear systems. Manuscript, December, 2008.
- [34] J. SHEN, L. HAN, AND J.S. PANG. Switching and stability properties of conewise linear systems. Revision under review, 2008. URL: <http://www.math.umbc.edu/~shenj/research/publications.html>
- [35] J. SHEN AND J.S. PANG. Linear complementarity systems with singleton properties: non-Zenoness. *Proceedings of 2007 American Control Conference*, pp.2769–2774, New York, 2007.
- [36] P. TAKAC. A short elementary proof of the Krein-Rutman Theorem. *Houston Journal of Mathematics*, Vol.20(1), pp.93–98, 1994.
- [37] A. TIWARI, J. FUNG, R. BHATTACHARYA, AND R.M. MURRAY. Polyhedral cone invariance applied to rendezvous of multiple agents. *Proceedings of the 43rd IEEE Conference on Decision and Control*, pp.165–170, Bahamas, 2004.
- [38] S. WIGGINS. *Introduction to Applied Dynamical Systems and Chaos*. Springer Verlag, 1990.
- [39] H. YAZAREL, AND G.J. PAPPAS. Geometric programming relaxations for linear system reachability. *Proceedings of 2004 American Control Conference*, pp.553–559, Boston, MA, 2004.
- [40] H. YAZAREL, S. PRAJNA, AND G.J. PAPPAS. S.O.S. for safety. *Proceedings of the 43rd IEEE Conference on Decision and Control*, pp.461–466, Bahamas, 2004.