



MIT Open Access Articles

Attention and Anticipation in Fast Visual-Inertial Navigation

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation	Carlone, Luca and Sertac Karaman. "Attention and Anticipation in Fast Visual-Inertial Navigation." Paper presented at the 2017 IEEE International Conference on Robotics and Automation (ICRA), Singapore, 29 May-3 June 2017, IEEE © 2017 The Author(s)
As Published	10.1109/TRO.2018.2872402
Publisher	Institute of Electrical and Electronics Engineers (IEEE)
Version	Original manuscript
Citable link	https://hdl.handle.net/1721.1/126592
Terms of Use	Creative Commons Attribution-Noncommercial-Share Alike
Detailed Terms	http://creativecommons.org/licenses/by-nc-sa/4.0/

Attention and Anticipation in Fast Visual-Inertial Navigation

Luca Carlone and Sertac Karaman

Abstract—We study a *Visual-Inertial Navigation* (VIN) problem in which a robot needs to estimate its state using an on-board camera and an inertial sensor, without any prior knowledge of the external environment. We consider the case in which the robot can allocate limited resources to VIN, due to tight computational constraints. Therefore, we answer the following question: under limited resources, what are the most relevant visual cues to maximize the performance of visual-inertial navigation? Our approach has four key ingredients. First, it is *task-driven*, in that the selection of the visual cues is guided by a metric quantifying the VIN performance. Second, it exploits the notion of *anticipation*, since it uses a simplified model for forward-simulation of robot dynamics, predicting the utility of a set of visual cues over a future time horizon. Third, it is *efficient and easy to implement*, since it leads to a greedy algorithm for the selection of the most relevant visual cues. Fourth, it provides *formal performance guarantees*: we leverage submodularity to prove that the greedy selection cannot be far from the optimal (combinatorial) selection. Simulations and real experiments on agile drones show that our approach ensures state-of-the-art VIN performance while maintaining a lean processing time. In the easy scenarios, our approach outperforms appearance-based feature selection in terms of localization errors. In the most challenging scenarios, it enables accurate visual-inertial navigation while appearance-based feature selection fails to track robot’s motion during aggressive maneuvers.

SUPPLEMENTARY MATERIAL

- Video: <https://www.youtube.com/watch?v=uMLXNRiVuyU>

I. INTRODUCTION

The human brain can extract conceptual information from an image in a time lapse as short as 13 ms [1]. One has proof of the human’s capability to seamlessly process large amount of sensory data in everyday tasks, including driving a car on a highway, or walking on a crowded street. In the cognitive science literature, there is agreement on the fact that efficiency in processing the large amount of data we are confronted with is due to our ability to prioritize some aspects of the visual scene, while ignoring others [2]. One can imagine that sensory inputs compete to have access to the limited computational resources of our brain. These resource constraints are dictated by the fixed amount of energy available to the brain as well as time constraints imposed by time-critical tasks. *Visual attention* is the cognitive process that allows humans to parse a large amount of visual data by selecting relevant information and filtering out irrelevant stimuli, so to maximize performance¹ under limited resources.

L. Carlone and S. Karaman are with the Laboratory for Information & Decision Systems (LIDS), Massachusetts Institute of Technology, Cambridge, MA, USA, {lcarlone,sertac}@mit.edu

¹This definition oversimplifies the attention mechanisms in humans. While the role of attention is to optimally allocate resources to maximize performance, it is known that some involuntary attention mechanisms can actually hinder the correct execution of a task [2].

Robots vs. humans. The astonishing progress in robotics and computer vision over the last three decades might induce us to ask: how far is robot perception from human performance? Let us approach this question by looking at the state of the art in visual processing for different tasks. Without any claim to be exhaustive, we consider few representative papers (sampled over the last 3 years) and we only look at timing performance. A state-of-the-art approach for object detection [3] detects objects in a scene in 22ms on a Titan X GPU. A high-performance approach for stereo reconstruction [4] builds a triangular mesh of a 3D scene in 10-100ms on a single CPU (at resolution 800×600). A state-of-the-art vision-based SLAM approach [5] requires around 400ms for local mapping and motion tracking and more than 1s for global map refinement (CPU, multiple cores). The reader may notice that for each task, in isolation, modern algorithms require more time than what a human needs to parse an entire scene. Arguably, while a merit of the robotics and computer vision communities has been to push performance in each task, we are quite far from a computational model in which all these tasks (pose estimation, geometry reconstruction, scene understanding) are concurrently executed in the blink of an eye.

Efficiency via general-purpose computing. One might argue that catching up with human efficiency is only a matter of time: according to Moore’s law, the available computational power grows at exponential rate, hence we only need to wait for more powerful computers. An analogous argument would suggest that using GPU rather than CPU would boost performance in some of the tasks mentioned above. By comparison with human performance, we realize that this argument is not completely accurate. While it is true that we can keep increasing the computational resources to meet given time constraints (i.e., enable faster processing of sensory data), the increase in computation implies an increase in energy consumption; for instance, a Titan X GPU has a nominal power consumption of 250W [6] while a Core i7 CPU has a power consumption as low as 11W [7]. On the other hand, human processing constantly deals with limited time and energy constraints, and is parsimonious in allocating only the resources that are necessary to accomplish its goals.

Efficiency via specialized computing. Another potential alternative to enable high-rate low-power perception and bridge the gap between human and robot perception is to design specialized hardware for machine perception. As extensively discussed in our previous work [8], algorithms and hardware co-design allows minimizing resource utilization by exploiting a tight-integration of algorithms and specialized hardware, and leveraging opportunities (e.g., pipelining, low-cost arithmetic) provided by ASICs (Application-Specific Integrated Circuits) and FPGAs (Field-Programmable Gate Arrays). While we

have shown that using specialized hardware for VIN leads to a reduction of the power consumption of 1-2 orders of magnitude (with comparable performance), three main observations motivate the present work. First, the development of specialized hardware for perception is an expensive and time-consuming process and the resulting hardware is difficult to upgrade. Second, rather than designing optimized hardware that can meet given performance requirements, it may be desirable to develop a framework that can systematically trade-off performance for computation, hence more flexibly adjusting to the available, possibly time-varying, computational resources and performance requirements. Third, extensive biological evidence suggests that efficient perception requires both specialized circuitry (e.g., visual perception in humans is carried out by highly specialized areas of the brain [9]) and a mechanism to prioritize stimuli (i.e., visual attention [2]).

Contribution. In this paper, we investigate how to speed-up computation (or, equivalently, reduce the computational effort) in visual navigation by prioritizing sensor data in a task-dependent fashion. In particular, we focus on a motion estimation task, *Visual-Inertial Navigation* (VIN), and consider the case in which, due to constraints on the on-board computation, a robot can only use a small number of visual features in the environment to support motion estimation. We then design a visual attention mechanism that selects a suitable set of visual features to maximize localization accuracy; our general framework is presented in Section III.

Our approach is task-driven: we consider a motion estimation task (VIN), and our approach selects features that maximize a task-dependent performance metric, that we present in Section III-A. Contrarily to the literature on visual feature selection, we believe that the utility of a feature is not an intrinsic property of the feature itself (e.g., appearance), but it rather stems from the intertwining of the environment and the observer state. Our approach seamlessly captures both visual saliency and the task-dependent utility of a set of features.

Our attention mechanism is predictive in nature: when deciding which feature is more useful, our approach performs fast forward-simulations of the state of the robot leading to a feature selection that is aware of the dynamics and the “intentions” of the robot. The forward simulation is based on a simplified model which we present in Section III-B.

Since the optimal allocation of the resources is a hard combinatorial problem, in Section IV we present a greedy algorithm for attention allocation. In the same section, we leverage recent results on submodularity to provide formal performance guarantees for the greedy algorithm. Section IV also reviews related techniques based on convex relaxations.

Section V provides an experimental evaluation of the proposed approach. The results confirm that our approach can boost performance in standard VIN pipelines and enables accurate navigation under agile motions and strict resource constraints. The proposed approach largely outperforms appearance-based feature selection methods, and drastically reduces the computational time required by the VIN back-end.

This paper extends the preliminary results presented in [10]. In particular, the discussion on convex relaxations for features

selection (Section IV-A and Section IV-B), the performance guarantees of Proposition 12, the simulation results of Section V-A, and the experimental evaluation on the 11 *EuRoC* datasets [11] are novel and have not been previously published.

II. RELATED WORK

This work intersects several lines of research across fields.

Attention and Saliency in Neuroscience and Psychology.

Attention is a central topic in human and animal vision research with more than 2500 papers published since the 1980s [2]. While a complete coverage is outside the scope of this work, we review few basic concepts, using the surveys of Carrasco [2], Borji and Itti [12], Scholl [13], and the work of Caduff and Timpf [14] as main references. Scholl [13] defines attention as the discrimination of sensory stimuli, and the allocation of limited resources to competing attentional demands. Carrasco [2] identifies three types of attention: *spatial*, *feature-based*, and *object-based*. Spatial attention prioritizes different locations of the scene by moving the eyes towards a specific location (*overt* attention) or by focusing on relevant locations within the field of view (*covert* attention). Feature-based attention prioritizes the detection of a specific feature (color, motion direction, orientation) independently on its location. Object-based attention prioritizes specific objects. In this work, we are mainly interested in covert spatial attention: which locations in the field of view are the most informative for navigation? Covert attention in humans is a combination of voluntary and involuntary mechanisms that guide the processing of visual stimuli at given locations in the scene [2]. Empirical evidence shows that attention is task-dependent in both primates and humans [14], [15]. Borji and Itti [12] explicitly capture this aspect by distinguishing bottom-up and top-down attention models; in the former the attention is captured by visual cues (stimulus-driven), while in the latter the attention is guided by the goal of the observer. Caduff and Timpf [14] study landmark saliency in human navigation and conclude that saliency stems from the intertwining of intrinsic property of a landmark (e.g., appearance) and the state of the observer (e.g., prior knowledge, observation pose). Another important aspect, that traces back to the *guided search theory* of Wolfe [16] and Spekreijse [17], is the distinction between pre-attentive and attentive visual processes. Pre-attentive processes handle *all* incoming sensory data in parallel; then, attentive processes only work on a filtered-out-version of the data, which the brain deems more relevant. General computational models for attention are reviewed in [12], including Bayesian models, graph-theoretic, and information-theoretic formulations.

Feature Selection in Robotics and Computer Vision.

The idea of enhancing performance in visual SLAM and visual odometry via active feature selection is not novel. Sim and Dudek [18] and Peretroukhin *et al.* [19] use training data to learn a model of the quality of the visual features. Each feature is mapped from a hand-crafted predictor space to a scalar weight that quantifies its usefulness for pose estimation; in [19] the weights are then used to rescale the measurement covariance of each observation. Ouerhani *et al.* [20] construct a topological map using attentional landmarks. Newman and

Ho [21] consider a robot equipped with camera and laser range finder and perform feature selection using an appearance-based notion of visual saliency. Sala *et al.* [22] use a co-visibility criterion to select good landmarks that are visible from multiple viewpoints. Siagian and Itti [23] investigate a bio-inspired attention model within Monte Carlo localization. Frintrop and Jensfelt [24] use an attention framework for landmark selection and active gaze control; feature selection is based on the VOCUS model [24], which includes a bottom-up attentional system (which computes saliency from the feature appearance), and can incorporate a top-down mechanism (which considers task performance). Active gaze control, instead, is obtained as the combination of three behaviors: landmark redetection, landmark tracking, and exploration of new areas. Hochdorfer and Schlegel [25] propose a landmark rating and selection mechanism based on area coverage to enable life-long mapping. Strasdat *et al.* [26] propose a reinforcement learning approach for landmark selection. Chli and Davison [27] and Handa *et al.* [28] use available priors to inform feature matching, hence reducing the computational cost. Jang *et al.* [29] propose an approach for landmark classification to improve accuracy in visual odometry; each feature class is used separately to estimate rotational and translational components of the ego-motion. Shi *et al.* [30] propose a feature selection technique to improve robustness of data association in SLAM. The notion of visual attention and saliency has been also investigated in the context of scene understanding. For instance, Oliva and Torralba [31] propose the notion of *Space Envelop* to obtain a coarse description of a scene which abstracts away unnecessary details, while Torralba *et al.* [32] propose an attention mechanism for natural search tasks, which combines bottom-up saliency and top-down aspects to identify the image region that is likely to attract the attention of human observer searching for a given object. Very recent work in computer vision use attention to reduce the computational burden in neural networks. Mnih *et al.* [33] reduce the processing of object detection and tracking with a recurrent neural network by introducing the notion of *glimpse*, which provides higher resolution in areas of interest within the image. Xu *et al.* [34] use visual attention to improve image content description. Cvišić and Petrović [35] speed up computation in stereo odometry by feature selection; the selection procedure is based on bucketing (which uniformly distributes the features across the image), and appearance-based ranking.

Our approach is loosely related to techniques for graph sparsification in which features are pruned a-posteriori from the SLAM factor graph to reduce computation; we refer the reader to the survey [36] for a review of these techniques.

The contributions that are most relevant to our proposal are the one of Davison [37], Lerner *et al.* [38], Mu *et al.* [39], Wu *et al.* [40], and Zhang and Vela [41]. The pioneering work of Davison [37] is one of the first papers to use information theoretic constructs to reason about visual features, and shares many of the motivations discussed in the present paper. Contrarily to the present paper, Davison [37] considers a model-based tracking problem in which the state of a

moving camera has to be estimated from observations of known features. In hindsight, we also provide a theoretical justification for the use of a greedy algorithm (similar to the one used in [37]) which we prove able to compute near-optimal solutions. Lerner *et al.* [38] study landmark selection in a localization problem with known landmarks; the robot has to choose a subset of landmarks to observe so to minimize the localization uncertainty. The optimal subset selection is formulated as a mixed-integer program and relaxed to an SDP. While the problem we solve in this paper is different (visual inertial odometry vs. localization with known map), an interesting aspect of [38] is the use of a *requirement matrix* that weights the state covariance and encodes task-dependent uncertainty constraints. Mu *et al.* [39] propose a two-stage approach to select a subset of landmarks to minimize the probability of collision and a subset of measurements to accurately localize those landmarks. Our approach shares the philosophy of task-driven measurement selection, but has three key differences. First, we use a simplified linear model for forward dynamics simulation: this is in the spirit of RANSAC, in that a simplified algebraic model is used to quickly filter out less relevant data. Second, we consider different performance metrics, going beyond the determinant criterion used in [39] and related work on graph sparsification. Third, we perform feature selection in a single stage and leverage submodularity to provide formal performance guarantees. Wu *et al.* [40] consider a multi camera system and split the feature selection process into a cascade of two resource-allocation problems: (i) how to allocate resources among the cameras, and (ii) how to select features in each camera, according to the allocated resources. The former problem is solved by taking simplifying assumptions on the distribution of the features, the latter is based on the heuristic feature selection scheme of [42]. Our paper attempts to formalize feature selection by leveraging the notion of submodularity. Zhang and Vela [41] perform feature selection using an observability score and provide suboptimality guarantees using submodularity. Our proposal is similar in spirit to [41] with few important differences. First, our approach is based on *anticipation*: the feature selector is aware of the intention (future motion) of the robot and selects the features accordingly. Second, we operate in a fixed-lag smoothing setup and investigate other performance metrics. Third, from the theoretical standpoint, we provide multiplicative suboptimality bounds and we prove conditions under which those bounds are non-vanishing.

Sensor Scheduling and Submodularity. Feature selection is deeply related to the problem of sensor scheduling in control theory. The most common setup for sensor scheduling is the case in which N sensors monitor a phenomenon of interest and one has to choose κ out of the N available sensors to maximize some information-collection metric; this setup is also known as *sensor selection* or *sensor placement*. The literature on sensor selection includes approaches based on convex relaxation [43], Bayesian optimal design [44], and submodular optimization [45]. The problem is shown to be NP-hard in [46]. Shamaiah *et al.* [47] leverage submodularity

and provide performance guarantees when optimizing the log-determinant of the estimation error covariance. A setup which is closer to the one in this paper is the case in which the sensed phenomenon is dynamic; in such case the sensor scheduling can be formulated in terms of the optimal selection of κ out of N possible measurements to be used in the update phase of a Kalman filter (KF). Vitus *et al.* [48] use a tree-search approach for sensor scheduling. Zhang *et al.* [49] proves that sensor scheduling within Kalman filtering is NP-hard and shows that the trace of the steady state prior and posterior KF covariances are not submodular, despite the fact that greedy algorithms are observed to work well in practice. Jawaid and Smith [50] provide counterexamples showing that in general the maximum eigenvalue and the trace of the covariance are not submodular. Tzoumas *et al.* [51] generalize the derivation of [50] to prove submodularity of the logdet over a fixed time horizon, under certain assumptions on the observation matrix. Summers *et al.* [52] show that several metrics based on the controllability and observability Gramians are submodular.

Visual-Inertial Navigation. As the combined use of the visual and vestibular system is key to human navigation, recent advances in visual-inertial navigation on mobile robots are enabling unprecedented performance in pose estimation in GPS-denied environments using commodity hardware. The literature on visual-inertial navigation is vast, with many contributions proposed over the last two years, including approaches based on filtering [53], [42], [54], [55], [56], [57], fixed-lag smoothing [58], [59], [60], [61], and full smoothing [62], [63], [64], [65], [66], [67], [68]. We refer the reader to [68] for a comprehensive review.

Notation. We use lowercase and uppercase bold letters to denote vectors (e.g. \mathbf{v}) and matrices (e.g. \mathbf{M}), respectively. Sets are denoted by sans script fonts (e.g. \mathbf{A}). Non-bold face letters are used for scalars and indices (e.g. j) and function names (e.g. $f(\cdot)$). The symbol $|\mathbf{A}|$ denotes the cardinality of \mathbf{A} . The identity matrix of size n is denoted with \mathbf{I}_n . An $m \times n$ zero matrix is denoted by $\mathbf{0}_{m \times n}$. $\mathbf{M} \succeq 0$ indicates that the matrix \mathbf{M} is positive semidefinite. The symbol $\|\cdot\|$ denotes the Euclidean norm for vectors and the spectral norm for matrices.

III. ATTENTION IN VISUAL-INERTIAL NAVIGATION

We design an attention mechanism that selects κ relevant visual features (e.g., Harris corners) from the current frame in order to maximize the performance of visual-inertial motion estimation. The κ features have to be selected out of N available features present in the camera image; the approach can deal with both monocular and stereo cameras (a stereo camera is treated as a rigid pair of monocular cameras).

We call \mathbf{F} the set of all available features (with $|\mathbf{F}| = N$). If we denote with $f(\cdot)$ our task-dependent performance metric (we formalize a suitable metric for VIN in Section III-A), we can state our feature selection problem as follows:

$$\max_{\mathbf{S} \subseteq \mathbf{F}} f(\mathbf{S}) \quad \text{subject to} \quad |\mathbf{S}| \leq \kappa \quad (1)$$

The problem looks for a subset of features \mathbf{S} , containing no more than κ features, which optimizes the task performance

$f(\cdot)$. This is a standard feature selection problem and has been used across multiple fields, including machine learning [69], robotics [39], and sensor networks [43]. Problem (1) is NP-hard [46] in general. In the rest of this paper we are interested in designing a suitable performance metric $f(\mathbf{S})$ for our VIN task, and provide fast approximation algorithms to solve (1).

We would like to design a performance metric $f(\cdot)$ that captures task-dependent requirements: in our case the metric has to quantify the uncertainty in the VIN motion estimation. Moreover, the metric should capture aspects already deemed relevant in related work. First, the metric has to reward the selection of the most distinctive features (in terms of appearance) since these are more likely to be re-observed in consecutive frames. Second, the metric has to reward features that remain within the field of view for a longer time. Therefore, *anticipation* is a key aspect: the metric has to be aware that under certain motion some of the features are more likely to remain in the field of view of the camera. Third, the metric has to reward features that are more informative to reduce uncertainty. In the following section we propose two performance metrics that seamlessly capture all these aspects.

A. Task-dependent Performance Metrics for VIN

Here we propose two metrics that quantify the accumulation of estimation errors over the horizon H , under the selection of a set of visual features \mathbf{S} . Assume that k is the time instant at which the features need to be selected. Let us call $\hat{\mathbf{x}}_k$ the (to-be-computed) state estimate of the robot at time k : we will be more precise about the variables included in $\hat{\mathbf{x}}_k$ in Section III-B1; for now the reader can think that $\hat{\mathbf{x}}_k$ contains the estimate for the pose and velocity of the robot at time k , as well as the IMU biases. We denote with $\hat{\mathbf{x}}_{k:k+H} \doteq [\hat{\mathbf{x}}_k \ \hat{\mathbf{x}}_{k+1} \ \dots \ \hat{\mathbf{x}}_{k+H}]$ the future state estimates within the horizon H . Moreover, we call $\mathbf{P}_{k:k+H}$ the covariance matrix of our estimate $\hat{\mathbf{x}}_{k:k+H}$, and $\mathbf{\Omega}_{k:k+H} \doteq \mathbf{P}_{k:k+H}^{-1}$ the corresponding information matrix. Two natural metrics to capture the accuracy of $\hat{\mathbf{x}}_{k:k+H}$ are described in the following.

Worst-case Estimation Error. The worst-case error variance is quantified by the largest eigenvalue $\lambda_{\max}(\mathbf{P}_{k:k+H})$ of the covariance matrix $\mathbf{P}_{k:k+H}$, see e.g., [43]. Calling $\lambda_{\min}(\mathbf{\Omega}_{k:k+H})$ the smallest eigenvalue of the information matrix $\mathbf{\Omega}_{k:k+H}$, it follows that $\lambda_{\max}(\mathbf{P}_{k:k+H}) = 1/\lambda_{\min}(\mathbf{\Omega}_{k:k+H})$, hence minimizing the worst-case error is the same as maximizing $\lambda_{\min}(\mathbf{\Omega}_{k:k+H})$. Note that the information matrix $\mathbf{\Omega}_{k:k+H}$ is function of the selected set of measurements \mathbf{S} , hence we write $\lambda_{\min}(\mathbf{\Omega}_{k:k+H}(\mathbf{S}))$.

Therefore our first metric (to be maximized) is:

$$f_{\lambda}(\mathbf{S}) = \lambda_{\min}(\mathbf{\Omega}_{k:k+H}(\mathbf{S})) = \lambda_{\min} \left(\bar{\mathbf{\Omega}}_{k:k+H} + \sum_{l \in \mathbf{S}} \Delta_l \right) \quad (2)$$

where on the right-hand-side, we exploited the additive structure of the information matrix, where $\bar{\mathbf{\Omega}}_{k:k+H}$ is the information matrix of the estimate when no features are selected (intuitively, this is the inverse of the covariance resulting from the IMU integration), while Δ_l is the information matrix

associated with the selection of the l -th feature. We will give an explicit expression to $\bar{\Omega}_{k:k+H}$ and Δ_l in Section III-B.

Volume and Mean Radius of the Confidence Ellipsoid. The ε -confidence ellipsoid is the ellipsoid that contains the estimation error with probability ε . Both the volume and the mean radius of the ε -confidence ellipsoid are proportional to the determinant of the covariance matrix. In particular, the volume V and the mean radius \bar{R} of an n -dimensional ellipsoid associated with the covariance $\mathbf{P}_{k:k+H}$ can be written as [43]:

$$V = \frac{(\alpha\pi)^{n/2}}{\Gamma(\frac{n}{2} + 1)} \det(\mathbf{P}_{k:k+H}^{\frac{1}{2}}), \quad \bar{R} = \sqrt{\alpha} \det(\mathbf{P}_{k:k+H})^{\frac{1}{2n}} \quad (3)$$

where α is the quantile of the χ^2 distribution with n degrees of freedom and upper tail probability of ε , $\Gamma(\cdot)$ is the Gamma function, and $\det(\cdot)$ is the determinant of a square matrix.

From (3) we note that to minimize the volume and the mean radius of the confidence ellipsoid we can equivalently minimize the determinant of the covariance. Moreover, since

$$\log \det(\mathbf{P}_{k:k+H}) = \log \det(\mathbf{\Omega}_{k:k+H}^{-1}) = -\log \det(\mathbf{\Omega}_{k:k+H})$$

then minimizing the size of the confidence ellipsoid is the same as maximizing the log-determinant of the information matrix, leading to our second performance metric:

$$f_{\det}(\mathbf{S}) = \log \det(\mathbf{\Omega}_{k:k+H}(\mathbf{S})) = \log \det \left(\bar{\Omega}_{k:k+H} + \sum_{l \in \mathbf{S}} \Delta_l \right) \quad (4)$$

where we again noted that the information matrix is function of the selected features and can be written in additive form.

Probabilistic Feature Tracks. The performance metrics described so far already capture some important aspects: they are task-dependent in that they both quantify the motion estimation performance; moreover, they are predictive, in the sense that they look at the result of selecting a set of features over a short (future) horizon. As we will see in Section III-B2, the model also captures the fact that longer feature tracks are more informative, therefore it implicitly rewards the selection of features that are co-visible across multiple frames.

The only aspect that is not yet modeled is the fact that, even when a feature is in the field of view of the camera, there is some chance that it will not be correctly tracked and the corresponding feature track will be lost. For instance, if the appearance of a feature is not distinctive enough, the feature track may be shorter than expected.

To model the probability that a feature track is lost, we introduce N Bernoulli random variables b_1, \dots, b_N . Each variable b_l represents the outcome of the tracking of feature l : if $b_l = 1$, then the feature is successfully tracked, otherwise, the feature track is lost. For each feature we assume $p_l = \text{Prob}(b_l = 1)$ to be given. In practice one can correlate the appearance of each feature to p_l , such that more distinctive features have higher probability of being tracked, or can learn the probabilities from data. Using the binary variables $\mathbf{b} \doteq \{b_1, \dots, b_N\}$, we write the information matrix at the end of the horizon as:

$$\mathbf{\Omega}_{k:k+H}(\mathbf{S}, \mathbf{b}) = \bar{\Omega}_{k:k+H} + \sum_{l \in \mathbf{S}} b_l \Delta_l \quad (5)$$

which has a clear interpretation: if the l -th feature is correctly tracked, then $b_l = 1$ and the corresponding information matrix Δ_l is added to $\bar{\Omega}_{k:k+H}$; on the other hand, if the feature tracks is lost, then $b_l = 0$ and the corresponding information content simply disappears from the sum in (5).

Since \mathbf{b} is a random vector, our information matrix is now a stochastic quantity $\mathbf{\Omega}_{k:k+H}(\mathbf{S}, \mathbf{b})$, hence we have to redefine our performance metrics to include the expectation over \mathbf{b} :

$$f(\mathbf{S}) = \mathbb{E}[f(\mathbf{\Omega}_{k:k+H}(\mathbf{S}, \mathbf{b}))] \quad (6)$$

where the function $f(\cdot)$ denotes either $f_{\lambda}(\cdot)$ or $f_{\det}(\cdot)$.

Computing the expectation (6) leads to a sum with a combinatorial number of terms, which is hard to even evaluate. To avoid the combinatorial explosion, we use Jensen's inequality, which produces an upper bound on the expectation of a *concave* function as follows:

$$\mathbb{E}[f(\mathbf{\Omega}_{k:k+H}(\mathbf{S}, \mathbf{b}))] \leq f(\mathbb{E}[\mathbf{\Omega}_{k:k+H}(\mathbf{S}, \mathbf{b})]) \quad (7)$$

Since both $f_{\lambda}(\cdot)$ and $f_{\det}(\cdot)$ are concave functions, the Jensen's inequality produces an upper bound for our expected cost. In the rest of this paper we maximize this bound, rather than the original cost. The advantage of doing so is that the right-hand-side of (7) can be efficiently computed as:

$$f(\mathbb{E}[\mathbf{\Omega}_{k:k+H}(\mathbf{S}, \mathbf{b})]) = f(\bar{\Omega}_{k:k+H} + \sum_{l \in \mathbf{S}} p_l \Delta_l) \quad (8)$$

where we used the definition (5), the fact that the expectation is a linear operator, and that $\mathbb{E}[b_l] = p_l$. Therefore, our performance metrics can be written explicitly as:

$$f_{\lambda}(\mathbf{S}) = \lambda_{\min} \left(\bar{\Omega}_{k:k+H} + \sum_{l \in \mathbf{S}} p_l \Delta_l \right) \quad (9)$$

$$f_{\det}(\mathbf{S}) = \log \det \left(\bar{\Omega}_{k:k+H} + \sum_{l \in \mathbf{S}} p_l \Delta_l \right)$$

which coincide with the deterministic counterparts (2), (4) when $p_l = 1, \forall l$. Interestingly, in (9) the probability that a feature is not tracked simply discounts the corresponding information content. Therefore, the approach considers features that are more likely to get lost as less informative, which is a desired behavior. We remark that the probabilities p_l only capture the distinctiveness of the visual features, while the geometric aspects (e.g., co-visibility) are captured in the matrices Δ_l , as described in Section III-B2. While the derivation so far is quite general and provides a feature selection mechanism for any feature-based SLAM system, in the following we focus on visual-inertial navigation and we provide explicit expressions for the matrices $\bar{\Omega}_{k:k+H}$ and Δ_l appearing in eq. (9).

B. Forward-simulation Model

The feature selection model proposed in Section III and the metrics in Section III-A require to predict the evolution of the information matrix over the horizon H . In the following we show how to forward-simulate the IMU and the camera; we note that we do not require to simulate actual IMU

measurements, but only need to predict the corresponding information matrix, which depends on the IMU noise statistics.

The forward-simulation model depends on the future motion of the robot (the IMU and vision models are function of the future poses of the robot); therefore, anticipation is a key element of our approach: the feature selection mechanism is aware of the immediate-future intentions of the robot and selects features accordingly. As we will see in the experiments, this enables a more clever selection of features during sharp turns and aggressive maneuvers. In practice, the future poses along the horizon can be computed from the current control actions; for instance, if the controller is planning over a receding horizon, one can get the future poses by integrating the dynamics of the vehicle. In this sense, our attention mechanism involves a tight integration of control and perception.

The algorithms for feature selection that we present in Section IV are generic and work for any positive definite $\mathbf{\Omega}_{k:k+H}$ and any positive semidefinite Δ_l . Therefore, the non-interested reader can safely skip this section, which provides explicit expressions for $\mathbf{\Omega}_{k:k+H}$ and Δ_l in the visual-inertial setup.

Before delving into the details of the IMU and vision model we remark a key design goal of our forward-simulation model: efficiency. The goal of an attention mechanism is to reduce the cognitive load later on in the processing pipeline; therefore, by design, it should not be computational demanding, as that would defeat its purpose. For this reasons, in this section we present a simplified VIN model which is designed to be efficient to compute, while capturing all the aspects of interest of a full visual-inertial estimation pipeline, e.g., [68].

1) *IMU Model and Priors*: Our simplified IMU model is based on a single assumption: the accumulation of the rotation error due to gyroscope integration over the time horizon is negligible. In other words, the relative rotation estimates predicted by the gyroscope are accurate. This assumption is realistic, even for inexpensive IMUs: the drift in rotation integration is typically small and negligible over the time horizon considered in our attention system (in our tests we consider a time horizon of 3s).

Assuming that the rotations are accurately known allows restricting the state to the robot position, linear velocity, and the accelerometer bias. Therefore, in the rest of this paper, the (unknown) state of the robot at time k is $\mathbf{x}_k \doteq [\mathbf{t}_k \ \mathbf{v}_k \ \mathbf{b}_k]$, where $\mathbf{t}_k \in \mathbb{R}^3$ is the 3D position of the robot, $\mathbf{v}_k \in \mathbb{R}^3$ is its velocity, and \mathbf{b}_k is the (time-varying) accelerometer bias. We also use the symbol \mathbf{R}_k to denote the attitude of the robot at time k : this is assumed to be known from gyroscope integration over the horizon H , hence it is not part of the state.

As in most VIN pipelines, we want to estimate the state of the robot at each frame². Therefore, the goal of this subsection, similarly to [68], is to reformulate a set of IMU measurements between two consecutive frames k and j as a single measurement that constrains \mathbf{x}_k and \mathbf{x}_j . Differently from [68], we show how to get a *linear* measurement model.

The on-board accelerometer measures the acceleration \mathbf{a}_k

of the sensor with respect to an inertial frame, and is affected by additive white noise $\boldsymbol{\eta}_k$ and a slowly varying sensor bias \mathbf{b}_k . Therefore, the measurement $\tilde{\mathbf{a}}_k \in \mathbb{R}^3$ acquired by the accelerometer at time k is modeled as [68]:

$$\tilde{\mathbf{a}}_k = \mathbf{R}_k^\top (\mathbf{a}_k - \mathbf{g}) + \mathbf{b}_k + \boldsymbol{\eta}_k, \quad (10)$$

where \mathbf{g} is the gravity vector, expressed in the inertial frame. To keep notation simple, we omit the reference frames in our notation, which follow closely the convention used in [68]: position \mathbf{t}_k and velocity \mathbf{v}_k are expressed in the (inertial) world frame,³ while the bias \mathbf{b}_k is expressed in the sensor frame.

Given position \mathbf{t}_k and velocity \mathbf{v}_k at time k , we can forward-integrate and obtain \mathbf{t}_j and \mathbf{v}_j at time $j > k$:

$$\begin{aligned} \mathbf{v}_j &= \mathbf{v}_k + \sum_{i=k}^{j-1} \mathbf{a}_i \delta \\ &\quad (\text{from (10) we know } \mathbf{a}_i = \mathbf{g} + \mathbf{R}_i(\tilde{\mathbf{a}}_i - \mathbf{b}_i - \boldsymbol{\eta}_i), \\ &\quad \text{and assuming constant bias between frames, } \mathbf{b}_i = \mathbf{b}_k) \\ &= \mathbf{v}_k + \mathbf{g} \delta_{kj} + \sum_{i=k}^{j-1} \mathbf{R}_i (\tilde{\mathbf{a}}_i - \mathbf{b}_k - \boldsymbol{\eta}_i) \delta \end{aligned} \quad (11)$$

$$\begin{aligned} \mathbf{t}_j &= \mathbf{t}_k + \sum_{i=k}^{j-1} (\mathbf{v}_i \delta + \frac{1}{2} \mathbf{a}_i \delta^2) \\ &\quad (\text{substituting } \mathbf{a}_i = \mathbf{g} + \mathbf{R}_i(\tilde{\mathbf{a}}_i - \mathbf{b}_k - \boldsymbol{\eta}_i)) \\ &= \mathbf{t}_k + \sum_{i=k}^{j-1} (\mathbf{v}_i \delta + \frac{1}{2} \mathbf{g} \delta^2 + \frac{1}{2} \mathbf{R}_i (\tilde{\mathbf{a}}_i - \mathbf{b}_k - \boldsymbol{\eta}_i) \delta^2) \\ &\quad (\text{substituting } \mathbf{v}_j \text{ from (11) with } j = i) \\ &= \mathbf{t}_k + \frac{1}{2} \mathbf{g} \delta_{kj}^2 + \sum_{i=k}^{j-1} \frac{1}{2} \mathbf{R}_i (\tilde{\mathbf{a}}_i - \mathbf{b}_k - \boldsymbol{\eta}_i) \delta^2 \\ &\quad + \sum_{i=k}^{j-1} (\mathbf{v}_k + \mathbf{g} \delta_{ki} + \sum_{h=k}^{i-1} \mathbf{R}_h (\tilde{\mathbf{a}}_h - \mathbf{b}_k - \boldsymbol{\eta}_h) \delta) \end{aligned} \quad (12)$$

where δ is the sampling time of the IMU, $\delta_{kj} \doteq \sum_{i=k}^{j-1} \delta$, and $\hat{\delta}_{kj}^2 \doteq \sum_{i=k}^{j-1} \delta^2$; as in [68], we assumed that the IMU bias remains constant between two frames. The evolution of the bias across frames can be modeled as a random walk:

$$\mathbf{b}_j = \mathbf{b}_k + \boldsymbol{\eta}_{kj}^b \quad (13)$$

where $\boldsymbol{\eta}_{kj}^b$ is a zero-mean random vector.

Noting that the state appears linearly in (11)-(13), it is easy to rewrite the three expressions together in matrix form:

$$\mathbf{z}_{kj}^{\text{IMU}} = \mathbf{A}_{kj} \mathbf{x}_{k:k+H} + \boldsymbol{\eta}_{kj}^{\text{IMU}} \quad (14)$$

where $\mathbf{z}_{kj}^{\text{IMU}} \in \mathbb{R}^9$ is a suitable vector,⁴ and $\boldsymbol{\eta}_{kj}^{\text{IMU}} \in \mathbb{R}^9$ is zero-mean random noise. We remark that while $\mathbf{z}_{kj}^{\text{IMU}}$ is function of the future IMU measurements, this vector is not actually used in our approach (what matters is \mathbf{A}_{kj} and the information matrix of $\boldsymbol{\eta}_{kj}^{\text{IMU}}$), hence we do not need to simulate future measurements. An explicit expression for the matrix $\mathbf{A}_{kj} \in \mathbb{R}^{9 \times 9(H+1)}$, the vector $\mathbf{z}_{kj}^{\text{IMU}}$, and the covariance of $\boldsymbol{\eta}_{kj}^{\text{IMU}}$ is given in Appendix A. The matrix \mathbf{A}_{kj} is a sparse block matrix with 9×9 blocks, which is all zeros, except the blocks corresponding to the state at times k and j .

³As usual when adopting MEMS inertial sensors, we assume that a local reference frame on Earth (our “world” frame) can be approximated as inertial, since the effects of the Earth rotation are negligible with respect to the measurement noise. For a more comprehensive discussion on reference frames for inertial navigation we refer the reader to [70, Chapter 2.2].

⁴The expression of $\mathbf{z}_{kj}^{\text{IMU}}$ is inconsequential for the following derivation, but the interested reader can find details and derivations in Appendix A.

² The derivation is identical for the case in which we associate a state to each keyframe, rather than each frame, as done in related work [68].

From linear estimation theory, we know that, using the IMU measurements (14) for all consecutive frames k, j in the horizon H , the information matrix of the optimal estimate of the state $\mathbf{x}_{k:k+H}$ given the IMU data is:

$$\bar{\Omega}_{k:k+H}^{\text{IMU}} = \sum_{k,j \in H} (\mathbf{A}_{kj}^T \Omega_{kj}^{\text{IMU}} \mathbf{A}_{kj}) \quad (15)$$

where H is the set of consecutive frames within the time horizon H , and $\Omega_{kj}^{\text{IMU}} \in \mathbb{R}^{9 \times 9}$ is the information matrix of the noise vector η_{kj}^{IMU} introduced in eq. (14).

While the IMU measurements constrain the states in the future horizon H , the predicted information matrix at time $k+H$ is also influenced by the initial information matrix at time k . This information matrix, referred to as $\bar{\Omega}_k^{\text{PRIOR}} \in \mathbb{R}^{9 \times 9}$, is computed and maintained by the VIN estimator,⁵ and can be understood as a prior on the state at time k . The presence of this information matrix results in an additional term in the expression of $\bar{\Omega}_{k:k+H}$, which is the information matrix of the state estimate before any vision measurement is selected, as per eq. (5). In particular, $\bar{\Omega}_{k:k+H}$ can be written as:

$$\bar{\Omega}_{k:k+H} = \bar{\Omega}_{k:k+H}^{\text{IMU}} + \bar{\Omega}_{k:k+H}^{\text{PRIOR}} \quad (16)$$

where $\bar{\Omega}_{k:k+H}^{\text{PRIOR}}$ is a matrix which is zero everywhere except the top-left 9×9 block which is equal to $\bar{\Omega}_k^{\text{PRIOR}}$, to reflect our prior on the state at time k . The matrix $\bar{\Omega}_{k:k+H}$ represents the information matrix of an optimal estimator of the state $\mathbf{x}_{k:k+H}$ before selecting visual measurements and follows the notation we have already introduced in (5). The presence of priors and measurements proceeds in full analogy with a standard fixed-lag smoother, while here we have the advantage of working with a linear model.

2) *Vision Model:* Also for the vision measurements, we are interested in designing a linear measurement model, which simplifies the actual (nonlinear) perspective projection model. To do so, we have to express a pixel measurement as a linear function of the unknown state that we want to estimate.

A (calibrated) pixel measurement of an external 3D point (or landmark) l identifies the 3D bearing of the landmark in the camera frame. Mathematically, if we call \mathbf{u}_{kl} the unit vector corresponding to the (calibrated) pixel observation of l from the robot pose at time k , \mathbf{u}_{kl} satisfies the following relation:

$$\mathbf{u}_{kl} \times ((\mathbf{R}_{\text{cam},k}^w)^T (\mathbf{p}_l - \mathbf{t}_{\text{cam},k}^w)) = \mathbf{0}_3 \quad (17)$$

where \times is the cross product between two vectors, \mathbf{p}_l is the 3D position of landmark l (in the world frame), $\mathbf{R}_{\text{cam},k}^w$ and $\mathbf{t}_{\text{cam},k}^w$ are the rotation and translation describing the camera pose at time k (w.r.t. the world frame). In words, the model (17) requires the observed point (transformed to the camera frame) to be collinear to the measured direction \mathbf{u}_{kl} , since the cross product measures the deviation from collinearity [71].

Now we note that for two vectors \mathbf{v}_1 and \mathbf{v}_2 , the cross product $\mathbf{v}_1 \times \mathbf{v}_2 = [\mathbf{v}_1]_{\times} \mathbf{v}_2$, where $[\mathbf{v}_1]_{\times}$ is the skew symmetric matrix built from \mathbf{v}_1 . Moreover, we note that the camera pose

w.r.t. the world frame, $(\mathbf{R}_{\text{cam},k}^w, \mathbf{t}_{\text{cam},k}^w)$, can be written as the composition of the IMU pose w.r.t. the world frame, $(\mathbf{R}_k, \mathbf{t}_k)$, and the relative pose of the camera w.r.t. the IMU, $(\mathbf{R}_{\text{cam}}^{\text{IMU}}, \mathbf{t}_{\text{cam}}^{\text{IMU}})$ (known from calibration). Using these two considerations, we rewrite (17) equivalently as:

$$[\mathbf{u}_{kl}]_{\times} ((\mathbf{R}_k \mathbf{R}_{\text{cam}}^{\text{IMU}})^T (\mathbf{p}_l - (\mathbf{t}_k + \mathbf{R}_k \mathbf{t}_{\text{cam}}^{\text{IMU}}))) = \mathbf{0}_3 \quad (18)$$

In presence of measurement noise, (18) becomes:

$$[\mathbf{u}_{kl}]_{\times} ((\mathbf{R}_k \mathbf{R}_{\text{cam}}^{\text{IMU}})^T (\mathbf{p}_l - (\mathbf{t}_k + \mathbf{R}_k \mathbf{t}_{\text{cam}}^{\text{IMU}}))) = \eta_{kl}^{\text{cam}} \quad (19)$$

where η_{kl}^{cam} is a zero-mean random noise with known covariance. Under the assumptions that rotations are known from gyroscope integration, the unknowns in model (19) are the robot position \mathbf{t}_k (which is part of our state vector $\mathbf{x}_{k:k+H}$) and the position of the observed 3D landmark \mathbf{p}_l . Rearranging the terms we obtain:

$$[\mathbf{u}_{kl}]_{\times} (\mathbf{R}_{\text{cam}}^{\text{IMU}})^T \mathbf{t}_{\text{cam}}^{\text{IMU}} = [\mathbf{u}_{kl}]_{\times} (\mathbf{R}_k \mathbf{R}_{\text{cam}}^{\text{IMU}})^T (\mathbf{t}_k - \mathbf{p}_l) + \eta_{kl}^{\text{cam}}$$

The only unknowns in the previous equation are \mathbf{t}_k and \mathbf{p}_l , hence the model is linear in the state and can be written in matrix form:

$$\mathbf{z}_{kl}^{\text{cam}} = \mathbf{F}_{kl} \mathbf{x}_{k:k+H} + \mathbf{E}_{kl} \mathbf{p}_l + \eta_{kl}^{\text{cam}} \quad (20)$$

with $\mathbf{z}_{kl}^{\text{cam}} \doteq [\mathbf{u}_{kl}]_{\times} (\mathbf{R}_{\text{cam}}^{\text{IMU}})^T \mathbf{t}_{\text{cam}}^{\text{IMU}}$, and for suitable matrices $\mathbf{F}_{kl} \in \mathbb{R}^{3 \times 9(H+1)}$ and $\mathbf{E}_{kl} \in \mathbb{R}^{3 \times 3}$. In order to be triangulated, a point has to be observed across multiple frames. Stacking the linear system (20) for each observation pose from which l is visible, we get a single linear system:

$$\mathbf{z}_l^{\text{cam}} = \mathbf{F}_l \mathbf{x}_{k:k+H} + \mathbf{E}_l \mathbf{p}_l + \eta_l^{\text{cam}} \quad (21)$$

where $\mathbf{z}_l^{\text{cam}} \in \mathbb{R}^{3n_\ell}$, $\mathbf{F}_l \in \mathbb{R}^{3n_\ell \times 9(H+1)}$, and $\mathbf{E}_l \in \mathbb{R}^{3n_\ell \times 3}$ are obtained by stacking (row-wise) $\mathbf{z}_{kl}^{\text{cam}}$, \mathbf{F}_{kl} , and \mathbf{E}_{kl} , respectively, for the n_ℓ frames in which l is visible. As for the IMU model, the expression of $\mathbf{z}_l^{\text{cam}}$ is inconsequential for our derivation, as it does not influence the future state covariance. On the other hand \mathbf{F}_l and \mathbf{E}_l depend on the future measurements \mathbf{u}_{kl} : for this reason, computing these matrices requires simulating pixel projections of \mathbf{p}_l for each frame in the horizon. When using a stereo camera, we have an estimate of \mathbf{p}_l hence we can easily project it to the future frames. In a monocular setup, we can guess the depth of new features from the existing features in the VIN back-end. We remark that when simulating pixel projections at future frames in the horizon, we also perform a visibility check (i.e., we use the camera projection model to assert whether a landmark projects within the image frame), hence restricting the set of future visual measurements according to the expected motion of the vehicle and the camera field of view. This aspect is crucial in making our feature selection approach “predictive” in nature. As discussed in Proposition 1 below, the resulting model seamlessly captures the intuitive fact that landmarks that remain visible and can be tracked across multiple frames are more informative.

Now we note that we cannot directly use the linear model (21) to estimate our state vector $\mathbf{x}_{k:k+H}$, since it contains the unknown position of landmark l . One way to

⁵ In our implementation, the VIN estimation is a fixed-lag smoother based on iSAM2 (described in Section V) and the information matrix at time k is obtained by marginalizing states in the smoother other than \mathbf{x}_k .

circumvent this problem is to include the 3D point in the state vector. This is undesirable for two reasons; first, including the landmarks as part of the state would largely increase the dimension of the state space (and hence of the matrices in (9)). Second, it may create undesirable behaviors of our performance metrics; for instance, the metrics might induce to select features that minimize the uncertainty of a far 3D point rather than focusing on the variables we are actually interested in (i.e., the state of the robot).

To avoid this undesirable effects, we analytically eliminate the 3D point from the estimation using the Schur complement trick [72]. We first write the information matrix of the joint state $[x_{k:k+H} \ p_l]$ from the linear measurements (21):

$$\Omega_{k:k+H}^{(l)} = \begin{bmatrix} F_l^\top F_l & F_l^\top E_l \\ E_l^\top F_l & E_l^\top E_l \end{bmatrix} \in \mathbb{R}^{(9H+12) \times (9H+12)} \quad (22)$$

where, for simplicity, we assumed the pixel measurement noise to be the identity matrix. Using the Schur complement trick we marginalize out the landmark l and obtain the information matrix of our state $x_{k:k+H}$ given the measurements (21):

$$\Delta_l = F_l^\top F_l - F_l^\top E_l (E_l^\top E_l)^{-1} E_l^\top F_l \in \mathbb{R}^{9(H+1) \times 9(H+1)} \quad (23)$$

Eq. (23) is the (additive) contribution to the information matrix of our state estimate due to the measurements of a single landmark l . This is the matrix that we already called Δ_l in (5). The matrix Δ_l is sparse, and its sparsity pattern is dictated by the co-visibility of landmark l across different frames [73]. It is worth noticing that $(E_l^\top E_l)^{-1}$ is the covariance of the estimate of the landmark position [73], and it is invertible as long as the landmark l can be triangulated. In our implementation, we only consider landmarks that can be triangulated (for which $(E_l^\top E_l)^{-1}$ invertible) as candidates for feature selection.

The following proposition formally proves that longer feature tracks lead to “larger” Δ_l hence the formulations (9) will encourage the selection of features with long tracks.

Proposition 1 (Long Feature Tracks). *Consider two landmarks l_1 and l_2 having identical predicted future pixel measurements till time k_1 (i.e., $z_{\tau l_1}^{cam} = z_{\tau l_2}^{cam}$) for $\tau = k, \dots, k_1$, but such that l_1 is only tracked till k_1 , while l_2 is tracked till a later frame k_2 (with $k < k_1 < k_2 \leq H$). Then $\Delta_{l_1} \preceq \Delta_{l_2}$.*

The proof of Proposition 1 is given in Appendix B. Proposition 1 ensures that features with long tracks carry a “larger” information content Δ_l hence contributing more to the objectives (9) (recall that for both choices of performance metrics it holds $f(\bar{\Omega}_{k:k+H} + \Delta_{l_1}) \leq f(\bar{\Omega}_{k:k+H} + \Delta_{l_2})$ whenever $\Delta_{l_1} \preceq \Delta_{l_2}$). Note, however, that Proposition 1 ensures that long feature tracks “dominate” short tracks only when their future measurements are identical. Therefore, in general, a heuristic selection based on feature track length may provide a suboptimal solution for (1) (intuitively, such choice would not account for the *geometry* of the features).

Remark 2 (Linear measurement models). *Sections III-B1 and III-B2 provide linear measurement models for inertial and visual measurements. In particular, we assumed that rotations are known over a short time horizon and this allowed us to*

obtain measurement models that are linear with respect to the unknown robot state. Within our framework, one might directly use linearized models of the nonlinear inertial and perspective models commonly used in VIN [68]. Our choice to design linear models has three motivations. First, we operate over a smaller state space (which does not include rotations and gyroscope biases), hence making matrix manipulations faster. Second, we avoid the actual computational cost of linearizing the nonlinear models. Third, thanks to the simplicity of the models, we enable a geometric understanding of our feature selection mechanisms (Section IV-D).

IV. ATTENTION ALLOCATION: ALGORITHMS AND PERFORMANCE GUARANTEES

In this section we discuss computational approaches to find a set of features that approximately solves the feature selection problem (1). It is known that finding the optimal subset S^* which solves (1) exactly is NP-hard [46], hence we cannot hope to find efficient algorithms to compute S^* in real-world problems.⁶ The solution we adopt in this paper is to design *approximation algorithms*, which are computationally efficient and provide performance guarantees (roughly speaking, produce a set which is not far from the optimal subset S^*). We remark that we are designing a *covert attention mechanism*: our algorithms only select a set of features that have to be retained and used for state estimation, while we do not attempt to actively control the motion of the camera.

In the following we present two classes of algorithms. The former is based on a convex relaxation of the original combinatorial problem (1). We do not claim the convex relaxations as original contributions, since they have been proposed multiple times in other contexts, see, e.g., [43]. We briefly review the convex relaxations and the corresponding performance guarantees in Section IV-A and Section IV-B. The second class of approximation algorithms includes greedy selection methods, and are discussed in greater details in Section IV-C. We provide performance guarantees for the greedy algorithms in Section IV-D.

A. Convex Relaxations

This section presents a convex-relaxation approach to compute an approximate solution for problem (1).

Using (9), we rewrite problem (1) explicitly as:

$$\max_{S \subseteq F} f(\bar{\Omega}_{k:k+H} + \sum_{l \in S} p_l \Delta_l) \quad \text{subject to } |S| \leq \kappa \quad (24)$$

where $f(\cdot)$ denotes either $f_\lambda(\cdot)$ or $f_{\det}(\cdot)$ (for the moment there is not need to distinguish the two metrics).

Introducing binary variables s_l , for $l = 1, \dots, N$, we rewrite (24) equivalently as:

$$\begin{aligned} & \max_{s_1, \dots, s_N} f(\bar{\Omega}_{k:k+H} + \sum_{l \in S} s_l p_l \Delta_l) \\ & \text{subject to } \sum_{l=1}^N s_l \leq \kappa, \quad s_l \in \{0, 1\} \quad \forall l \in \{1, \dots, N\} \end{aligned} \quad (25)$$

⁶ In typical real-world problems, the set of available visual feature is larger than 200, and we are asked to select 10–100 features, depending on on-board resources. In those instances, the cost of a brute force search is prohibitive.

Problem (25) is a binary optimization problem. While problem (25) would return the optimal subset \mathbf{S}^* , it is still NP-hard to solve, due to the constraint that s_l have to be binary.

Problem (25) admits a simple convex relaxation:

$$f_{\text{cvx}}^* = \max_{s_1, \dots, s_N} f(\bar{\Omega}_{k:k+H} + \sum_{l \in \mathbf{S}} s_l p_l \Delta_l) \quad (26)$$

subject to $\sum_{l=1}^N s_l \leq \kappa$, $s_l \in [0, 1] \forall l \in \{1, \dots, N\}$

where the binary constraint $s_l \in \{0, 1\}$ is replaced by the convex constraint $s_l \in [0, 1]$. Convexity of problem (26) follows from the fact that we maximize a concave cost under linear inequality constraints.⁷

This convex relaxation has been proposed multiple times in other contexts (see, e.g., [43]). The solution s_1^*, \dots, s_N^* of (26) is not binary in general and a rounding procedure is needed to distinguish the features that have to be discarded ($s_l = 0$) from the features that have to be selected ($s_l = 1$). A common rounding procedure is to simply select the κ features with the largest s_l , while randomized rounding procedures have also been considered [38]. We use the former, and we call \mathbf{S}° the set including the indices of the κ features with the largest s_l^* , where s_1^*, \dots, s_N^* is the optimal solution of (26).

B. Performance Guarantees for the Convex Relaxations

The convex relaxation (26) has been observed to work well in practice, although there is no clear (a-priori) performance guarantee on the quality of the set \mathbf{S}° .

Let us call f_{cvx}^* the optimal objective of the relaxed problem (26), $f(\mathbf{S}^\circ)$ the objective attained by the rounded solution, and $f(\mathbf{S}^*)$ the optimal solution of the original NP-hard problem (25). Then, one can easily obtain a-posteriori performance bounds by observing that:

$$f(\mathbf{S}^\circ) \leq f(\mathbf{S}^*) \leq f_{\text{cvx}}^* \quad (27)$$

where the first inequality follows from optimality of \mathbf{S}^* (any subset of κ features has cost at most $f(\mathbf{S}^*)$), while the latter from the fact that (26) is a relaxation of the original problem.

The chain of inequality (27) suggests a simple (a-posteriori) performance bound for the quality of the set produced by the convex relaxation (26):

$$f(\mathbf{S}^*) - f(\mathbf{S}^\circ) \leq f_{\text{cvx}}^* - f(\mathbf{S}^\circ) \quad (28)$$

i.e., the suboptimality gap $f(\mathbf{S}^*) - f(\mathbf{S}^\circ)$ of the subset \mathbf{S}° is bounded by the difference $f_{\text{cvx}}^* - f(\mathbf{S}^\circ)$, which can be computed (a posteriori) after solving (26).

C. Greedy Algorithms and Lazy Evaluation

This section presents a second approach to approximately solve problem (1). Contrarily to the convex relaxation of Section IV-A, here we consider a greedy algorithm that selects κ features that (approximately) maximize the cost $f(\cdot)$.

The algorithm starts with an empty set $\mathbf{S}^\#$ and performs κ iterations. At each iteration, it adds the feature that, if added to $\mathbf{S}^\#$, induces the largest increase in the cost function. The pseudocode of the algorithm is given in Algorithm 1.

⁷Both the smallest eigenvalue and the log-determinant of a positive definite matrix are concave functions [74] of the matrix entries.

Algorithm 1: Greedy algorithm with lazy evaluation

```

1 Input:  $\bar{\Omega}_{k:k+H}$ ,  $\Delta_l$ , for  $l = 1, \dots, N$ , and  $\kappa$ ;
2 Output: feature subset  $\mathbf{S}^\#$ ;
3  $\mathbf{S}^\# = \emptyset$ ;
4 for  $i = 1, \dots, \kappa$  do
5   % Compute upper bound for  $f(\mathbf{S}^\# \cup \{l\})$ ,  $l = 1, \dots, N$ 
6    $[U_1, \dots, U_N] = \text{upperBounds}(\bar{\Omega}_{k:k+H}, \Delta_1, \dots, \Delta_N)$ ;
7   % Sort features using upper bound
8    $\mathbf{F}^\downarrow = \text{sort}(U_1, \dots, U_N)$ ;
9   % Initialize best feature
10   $f_{\text{max}} = -1$ ;  $l_{\text{max}} = -1$ ;
11  for  $l \in \mathbf{F}^\downarrow$  do
12    if  $U_l < f_{\text{max}}$  then
13      break;
14    end
15    if  $f(\mathbf{S}^\# \cup l) > f_{\text{max}}$  then
16       $f_{\text{max}} = f(\mathbf{S}^\# \cup l)$ ;  $l_{\text{max}} = l$ ;
17    end
18  end
19   $\mathbf{S}^\# = \mathbf{S}^\# \cup l_{\text{max}}$ ;
20 end

```

In line 3 the algorithm starts with an empty set. The “for” loop in line 4 iterates κ times: at each time the best feature is added to the subset $\mathbf{S}^\#$ (line 19). The role of the “for” loop in line 11 is to compute the feature that induces the maximum increase in the cost (lines 15-17). The remaining lines provide a lazy evaluation mechanism. For each feature l we compute an upper bound on the cost $f(\mathbf{S}^\# \cup \{l\})$ (line 6). The features are sorted (in descending order) according to this upper bound (line 8). The advantage of this is that by comparing the current best feature with this upper bound (line 12) we can avoid checking features that are guaranteed to attain a smaller cost.

Clearly, the lazy evaluation is advantageous if the upper bound is faster to compute than the actual cost. The following propositions provide two useful (and computationally cheap) upper bounds for our cost functions.

Proposition 3 (Upper bounds for log det: Hadamard’s inequality, Thm 7.8.1 [75]). *For a positive definite matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$ with diagonal elements M_{ii} , it holds:*

$$\det(\mathbf{M}) \leq \prod_{i=1}^n M_{ii} \Leftrightarrow \log \det(\mathbf{M}) \leq \sum_{i=1}^n \log M_{ii} \quad (29)$$

Proposition 4 (Eigenvalue Perturbation Bound [76]). *Given Hermitian matrices $\mathbf{M}, \Delta \in \mathbb{R}^{n \times n}$, and denoting with $\lambda_i(\mathbf{M})$ the i -th eigenvalue of \mathbf{M} , the following inequalities hold:*

$$|\lambda_i(\mathbf{M} + \Delta) - \lambda_i(\mathbf{M})| \leq \|\Delta\| \quad (30)$$

$$\min_j |\lambda_i(\mathbf{M}) - \lambda_j(\mathbf{M} + \Delta)| \leq \|\Delta \mathbf{v}_i\| \quad (31)$$

where \mathbf{v}_i is the eigenvector of \mathbf{M} associated to $\lambda_i(\mathbf{M})$.

Eq. (30) is a restatement of the classical Weyl inequality, while (31) is a tighter bound from Ipsen and Nadler [76]. To clarify how the bounds in Proposition 4 provide us with an upper bound for λ_{\min} , we prove the following result.

Corollary 5 (Upper bounds for λ_{\min}). *Given two symmetric and positive semidefinite matrices $\mathbf{M}, \Delta \in \mathbb{R}^{n \times n}$ the following inequality holds:*

$$\lambda_{\min}(\mathbf{M} + \Delta) \leq \lambda_{\min}(\mathbf{M}) + \|\Delta \mathbf{v}_{\min}\| \quad (32)$$

where \mathbf{v}_{\min} is the eigenvector of \mathbf{M} associated to the smallest eigenvalue $\lambda_{\min}(\mathbf{M})$.

The proof of Corollary 5 is given in Appendix C. While Algorithm 1 highlights the simplicity of the greedy algorithm, it is unclear whether this algorithm produces good subsets of features. We tackle this question in the next section.

D. Performance Guarantees for the Greedy Algorithm

This section shows that the greedy algorithm (Algorithm 1) admits provable sub-optimality bounds. These bounds guarantee that the greedy selection cannot perform much worse than the optimal strategy. The section tackles separately the two metrics presented in Section III-A, since the corresponding performance guarantees are fundamentally different.

Our results are based on the recent literature on submodularity and submodular maximization. Before delving in the guarantees for each metric, we provide few preliminary definitions, which can be safely skipped by the expert reader.

Definition 6 (Normalized and Monotone Set Function [77]). *A set function $f : 2^F \rightarrow \mathbb{R}$ is said to be normalized if $f(\emptyset) = 0$; $f(\mathcal{S})$ is said to be monotone (non-decreasing) if for any subsets $A \subseteq B \subseteq F$, it holds $f(A) \leq f(B)$.*

Definition 7 (Submodularity [77]). *A set function $f : 2^F \rightarrow \mathbb{R}$ is submodular if, for any subsets $A \subseteq B \subseteq F$, and for any element $e \in F \setminus B$, it holds that:*

$$f(A \cup \{e\}) - f(A) \geq f(B \cup \{e\}) - f(B) \quad (33)$$

Submodularity formalizes the notion of diminishing returns: adding a measurement to a small set of measurement is more advantageous than adding it to a large set. Our interest towards submodularity is motivated by the following result.

Proposition 8 (Near-optimal submodular maximization [77]). *Given a normalized, monotone, submodular set function $f : 2^F \rightarrow \mathbb{R}$, and calling \mathcal{S}^* the optimal solution of the maximization problem (1), then the set $\mathcal{S}^\#$, computed by the greedy Algorithm 1, is such that:*

$$f(\mathcal{S}^\#) \geq (1 - 1/e)f(\mathcal{S}^*) \approx 0.63f(\mathcal{S}^*) \quad (34)$$

This bound ensures us that the worst-case performance of the greedy algorithm cannot be far from the optimum. In the following we tailor this result to our feature selection problem.

1) *Sub-optimality guarantees for $\log \det$:* It is possible to show that $\log \det$ is submodular with respect to the set of measurements used for estimation. This result and the corresponding performance guarantees are formalized as follows.

Proposition 9 (Submodularity of $\log \det$ [47]). *The set function $f_{\det}(\mathcal{S})$ defined in (4) is monotone and submodular.*

Moreover, the greedy algorithm applied to (1) using $f_{\det}(\mathcal{S})$ as objective enjoys the following performance guarantees:

$$f_{\det}(\mathcal{S}) \geq (1 - 1/e)f_{\det}(\mathcal{S}^*) + \frac{f_{\det}(\emptyset)}{e} \quad (35)$$

The result is proven in [47] and has been later rectified to account for the need of normalized functions in [50]. The extra term $\frac{f_{\det}(\emptyset)}{e}$ in (35) indeed follows from the application of Proposition 8 to the normalized function $f_{\det}(\mathcal{S}) - f_{\det}(\emptyset)$.

2) *Sub-optimality guarantees for λ_{\min} :* Currently, no result is readily available to bound the suboptimality gap of the greedy algorithm applied to the maximization of the smallest eigenvalue of the information matrix (or equivalently minimizing the largest eigenvalue of the covariance). Indeed, related work provides counterexamples, showing that this metric is not submodular in general, while the greedy algorithm is observed to perform well in practice [50]. In this section we provide a first result showing that, despite the fact that $f_{\lambda}(\mathcal{S})$ fails to be submodular, it is not far from a submodular function. This notion is made more precise in the following.

Definition 10 (Submodularity ratio [69], [78]). *The submodularity ratio of a non-negative set function $f(\cdot)$ with respect to a set \mathcal{S} and an integer $\kappa \geq 1$ is defined as:*

$$\gamma_{\mathcal{S}} \doteq \min_{\substack{L \subseteq \mathcal{S}, \\ E: |E| \leq \kappa, E \cap L = \emptyset}} \frac{\sum_{e \in E} (f(L \cup \{e\}) - f(L))}{f(L \cup E) - f(L)} \quad (36)$$

It is possible to show that if $\gamma_{\mathcal{S}} \geq 1$, then the function $f(\cdot)$ is submodular. However, in this context we are interested in the submodularity ratio, since it enables a less restrictive sub-optimality bound, as described in the following proposition.

Proposition 11 (Approximate submodular maximization [69], [78]). *Let $f(\cdot)$ be a non-negative monotone set function and let \mathcal{S}^* be the optimal solution of the maximization problem (1), then the set $\mathcal{S}^\#$, computed by the greedy Algorithm 1 is such that:*

$$f(\mathcal{S}^\#) \geq (1 - e^{-\gamma_{\mathcal{S}^\#}})f(\mathcal{S}^*) \quad (37)$$

where $\gamma_{\mathcal{S}^\#}$ is the submodularity ratio of $f(\cdot)$ with respect to $\mathcal{S}^\#$ and $\kappa = |\mathcal{S}^\#|$.

Proposition 11 provides a multiplicative suboptimality bound whenever $\gamma_{\mathcal{S}^\#} > 0$. In the following we show that this is generally the case when maximizing the smallest eigenvalue.

Proposition 12 (Non-vanishing Submodularity ratio of λ_{\min}). *Call $\mathcal{S}^\#$ the set returned by the greedy algorithm maximizing λ_{\min} . For any set $L \subseteq \mathcal{S}^\#$ call $\bar{\mu}$ the eigenvector corresponding to the smallest eigenvalue of the matrix $\bar{\Omega}_{k:k+H} + \sum_{l \in L} \Delta_l$. Moreover call $\bar{\mu}_0, \bar{\mu}_2, \dots, \bar{\mu}_H \in \mathbb{R}^3$, the subvectors of $\bar{\mu}$ corresponding the robot positions at time $k, \dots, k+H$. Then the submodularity ratio of λ_{\min} is bounded away from zero if $\bar{\mu}_i \neq \bar{\mu}_j$, for some i, j .*

The proof of Proposition 12 is given in Appendix D. In words, Proposition 12 states that the submodularity ratio does not vanish as long as the directions of largest uncertainty

change along the future horizon. The following corollary is a straightforward consequence of Proposition 12.

Corollary 13 (Approximate submodularity of λ_{\min}). *The set function $f_\lambda(\mathbf{S})$ defined in (9) is monotone. Moreover, under the assumptions of Proposition 12, the greedy algorithm applied to (1) using $f_\lambda(\mathbf{S})$ as objective enjoys the guarantees of Proposition 11 for a nonzero $\gamma_{\mathbf{S}^\#}$.*

Proof: Monotonicity follows from the Weyl inequality [76]. The guarantees of the greedy algorithm follow from Proposition 11 and Proposition 12. ■

Corollary 13 guarantees that the approximation bound of Proposition 11 does not vanish, hence the greedy algorithm always approximate the optimal solution up to a constant-factor. Empirical evidence, shown in Section V, confirms that the greedy algorithm applied to the maximization of $f_\lambda(\mathbf{S})$ has excellent performance, producing near-optimal results in all test instances.

Remark 14 (Geometric Intuition Behind Greedy Selection with λ_{\min}). *Our linear model enables a deeper understanding of the geometry behind the greedy selection. The greedy selection rewards features l with large objective $\lambda_{\min}(\bar{\mathbf{\Omega}}_{k:k+H} + \Delta_l)$ or, equivalently, large marginal gain $\lambda_{\min}(\bar{\mathbf{\Omega}}_{k:k+H} + \Delta_l) - \lambda_{\min}(\bar{\mathbf{\Omega}}_{k:k+H})$. The following chain of relations provides a geometric understanding of which features induce a large marginal gain:*

$$\begin{aligned} & \lambda_{\min}(\bar{\mathbf{\Omega}}_{k:k+H} + \Delta_l) - \lambda_{\min}(\bar{\mathbf{\Omega}}_{k:k+H}) \\ & \text{(from Rayleigh quotient)} \\ & = \min_{\|\boldsymbol{\nu}\|=1} \boldsymbol{\nu}^\top (\bar{\mathbf{\Omega}}_{k:k+H} + \Delta_l) \boldsymbol{\nu} - \min_{\|\boldsymbol{\mu}\|=1} \boldsymbol{\mu}^\top (\bar{\mathbf{\Omega}}_{k:k+H}) \boldsymbol{\mu} \\ & \text{(calling } \bar{\boldsymbol{\mu}} \text{ the minimizer of the second summand)} \\ & = \min_{\|\boldsymbol{\nu}\|=1} \boldsymbol{\nu}^\top (\bar{\mathbf{\Omega}}_{k:k+H} + \Delta_l) \boldsymbol{\nu} - \bar{\boldsymbol{\mu}}^\top (\bar{\mathbf{\Omega}}_{k:k+H}) \bar{\boldsymbol{\mu}} \\ & \text{(substituting a suboptimal solution } \bar{\boldsymbol{\mu}} \text{ in the first summand)} \\ & \leq \bar{\boldsymbol{\mu}}^\top (\bar{\mathbf{\Omega}}_{k:k+H} + \Delta_l) \bar{\boldsymbol{\mu}} - \bar{\boldsymbol{\mu}}^\top (\bar{\mathbf{\Omega}}_{k:k+H}) \bar{\boldsymbol{\mu}} \\ & \text{(simplifying and substituting the expression of } \Delta_l) \\ & = \bar{\boldsymbol{\mu}}^\top \Delta_l \bar{\boldsymbol{\mu}} = \bar{\boldsymbol{\mu}}^\top \mathbf{F}_l^\top (\mathbf{I} - \mathbf{E}_l (\mathbf{E}_l^\top \mathbf{E}_l)^{-1} \mathbf{E}_l^\top) \mathbf{F}_l \bar{\boldsymbol{\mu}} \\ & \text{(defining the idempotent matrix } \mathbf{Q} \doteq (\mathbf{I} - \mathbf{E}_l (\mathbf{E}_l^\top \mathbf{E}_l)^{-1} \mathbf{E}_l^\top)) \\ & = \bar{\boldsymbol{\mu}}^\top \mathbf{F}_l^\top \mathbf{Q} \mathbf{F}_l \bar{\boldsymbol{\mu}} = \bar{\boldsymbol{\mu}}^\top \mathbf{F}_l^\top \mathbf{Q} \mathbf{Q}^\top \mathbf{F}_l \bar{\boldsymbol{\mu}} = \|\mathbf{Q} \mathbf{F}_l \bar{\boldsymbol{\mu}}\|^2 \\ & \text{(using the triangle inequality and substituting } \mathbf{F}_l) \\ & \leq \|\mathbf{Q}\|^2 \|\mathbf{F}_l \bar{\boldsymbol{\mu}}\|^2 = \|\mathbf{F}_l \bar{\boldsymbol{\mu}}\|^2 = \sum_{k=0}^H \|\mathbf{u}_{kl}\| \times (\mathbf{R}_{cam,k}^w)^\top \bar{\boldsymbol{\mu}}_k \|^2 \end{aligned}$$

where $\bar{\boldsymbol{\mu}}_k$ is the subvector of $\bar{\boldsymbol{\mu}}$ at the entries corresponding to the robot position coordinates at time k . Intuitively, the inequalities reveal that the marginal gain is small when $\|\mathbf{u}_{kl}\| \times (\mathbf{R}_{cam,k}^w)^\top \bar{\boldsymbol{\mu}}_k$ is small, i.e., when we pick landmark observations where the measured bearing \mathbf{u}_{kl} is nearly parallel to the directions of large uncertainty $\bar{\boldsymbol{\mu}}_k$, transformed in the camera frame by the rotation $(\mathbf{R}_{cam,k}^w)^\top$. For instance, if we have large uncertainty in the forward direction, it is not convenient to use features in front of the robot (i.e., with bearing parallel to the direction of largest uncertainty); accordingly, the greedy approach would select features in the periphery of the image, which intuitively provide a better way to reduce uncertainty.

V. EXPERIMENTS

This section provides three sets of experimental results. The first set of tests, in Section V-A, shows that the greedy algo-

rithm attains near-optimal solutions in solving problem (1), while being faster than general purpose solvers for the convex relaxations discussed in Section IV-A. The second set of tests, in Section V-B, evaluates our c++ pipeline in realistic simulations, showing that our feature selection techniques boost VIN performance; the same section also shows the advantage of using our lazy evaluation. The third set of tests, in Section V-C, evaluates our approach on real data collected by an agile micro aerial vehicle.

A. Assessment of the Greedy Algorithms for Feature Selection

This section answers the following question: *how good is the greedy Algorithm 1 to (approximately) solve the combinatorial optimization problem (1)?* In particular, we show that the greedy algorithm finds a near-optimal solution of (1), for both choices of the cost function (9); we also show that the convex relaxation approach of Section IV-A finds near-optimal solutions, while being more computationally expensive.

Testing setup. To generate random instances of problem (1), we consider a robot moving along a straight line at a constant speed of 2m/s. The robot is equipped with an IMU with sampling period $\delta = 0.01$ s; we choose the accelerometer noise density equal to $0.02\text{m}/(\text{s}^2\sqrt{\text{Hz}})$, and the accelerometer bias continuous-time noise density to be $0.03\text{m}/(\text{s}^3\sqrt{\text{Hz}})$. We also simulate an on-board monocular camera, which measures 3D points randomly scattered in the environment, at a (key)frame rate of 0.5s. The robot has to select a set of κ features out of N available visual measurements. We assume that at the time of feature selection, the position covariance of the robot is $10^{-2} \cdot \mathbf{I}_3$, while its velocity and accelerometer bias covariances are $10^{-2} \cdot \mathbf{I}_3$ and $10^{-4} \cdot \mathbf{I}_3$, respectively. Using this information, we build the matrix $\bar{\mathbf{\Omega}}_{k:k+H}$, using a prediction horizon of 2.5s. Moreover, from the available feature measurements, we build the matrices Δ_l ; in these tests we assume $p_l = 1$, i.e., we disregard appearance during feature selection.

Techniques and evaluation metrics. We compare two approaches to solve (1): the greedy algorithm of Algorithm 1 and the convex relaxation approach (26). We implemented the convex relaxation using CVX/MOSEK as parser/solver for (26), and then we computed the rounded solution as described in Section IV-A. For the evaluation in this section, we implemented both the greedy algorithm and the convex relaxation in Matlab. We evaluate these approaches for each choice of the objective functions f_λ and f_{\det} defined in (9). Ideally, for each technique, we should compare the objective attained by the techniques, versus the optimal objective. Unfortunately, the optimal objective is hard to compute and a brute-force approach is prohibitively slow, even for relatively small problem instances.⁸ Luckily, the convex relaxation (26) also produces an upper bound on the optimal cost of (1) (c.f. eq. (27)), hence we can use this upper bound to understand how far are the greedy and the rounded solution of (26) from optimality.

Results. We consider problems of increasing sizes in which we are given N features and we have to select half of them

⁸Even in a small instance in which we are required to select 50 out of 100 available visual measurements, a brute-force approach would need to evaluate around 10^{29} possible sets.

($\kappa = N/2$) to maximize the objective in (1). For each N , we compute statistics over 50 Monte Carlo.

Fig. 1(a) shows the smallest eigenvalue objective f_λ attained by the different techniques for increasing number of features N . Besides the greedy, the rounded convex relaxation (label: rounded), and the relaxed objective (label: relaxed), we show the objective attained by picking a random subset of κ features (label: random). We are solving a maximization problem hence the larger the objective the better. Fig. 1(a) shows that in all tested instances, greedy and rounded match the upper bound relaxed (the three lines are practically indistinguishable), hence they both produce optimal solutions (c.f. eq. (28)). The resulting solution is far better than random. This result is somehow surprising, since the smallest eigenvalue is not submodular in general, and the greedy algorithm enjoys weaker performance guarantees (Corollary 13). However, this observation is in agreement with related work in other fields, e.g., [49]. While both greedy and rounded return good solutions, solving the convex problem (26) is usually more expensive than computing the greedy solution: the CPU time of our greedy algorithm in Matlab (without lazy evaluation) is around 0.4s (for $N = 50$), while CVX requires around 0.8s.

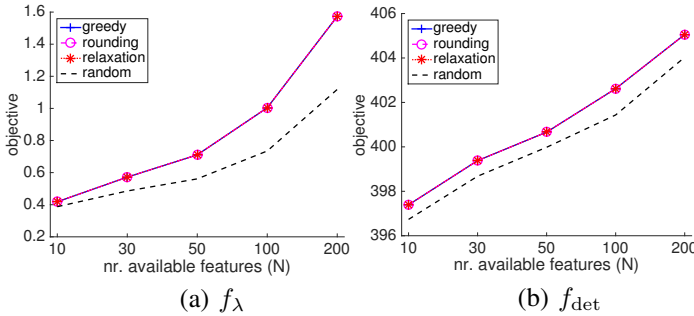


Fig. 1. Techniques to approximately solve problem (1) for (a) the smallest eigenvalue objective f_λ , and (b) the log-determinant objective f_{\det} . The figure reports the objective attained by the greedy algorithm (greedy), the rounded solution (rounded), and a random selection (random). The upper bound relaxed, attained by the convex problem (26) (before rounding), is shown for comparison.

Analogous considerations hold for the objective f_{\det} . Fig. 1(b) shows the log-determinant attained by the different techniques, for increasing number of available features N ; also in this case the algorithms have to select $\kappa = N/2$ features. As in the previous tests, greedy and rounded attain the optimal solution in all test instances, matching the upper bound relaxed, and performing remarkably better than a random choice. Regarding the CPU time, our Matlab implementation of the greedy algorithm to optimize f_{\det} takes around 0.1s (for $N = 50$), while CVX requires more than 1min to solve (26).⁹ We remark that while CVX is a state-of-the-art general-purpose solver for convex programming, our analysis does not rule out the possibility of designing fast specialized solvers, e.g., [79]; such an attempt, however, is outside the scope of this paper. Since the greedy algorithms are as accurate as the convex

⁹CVX uses a successive approximation method to maximize the log-det objective, which is known to be fairly slow.

relaxation techniques, while being faster than general-purpose convex solvers, in the following we focus on the former.

B. Importance of Feature Selection in VIN

This section answers the following question: *does the feature selection resulting by solving (1) lead to performance improvements in VIN?* In the following we show that the proposed feature selection approach boosts VIN performance in realistic Monte Carlo simulations.

Testing setup. We adopt the benchmarking problem of [68] and pictured in Fig. 2(a) as testing setup. We simulate a robot that follows a circular trajectory with a sinusoidal vertical motion. The total length of the trajectory is 120m. The on-board camera has a focal length of 315 pixels and runs at a rate of 2.5Hz (simulating keyframes). Simulated acceleration and gyroscope measurements are obtained as in [68].

Implementation and evaluation metrics. In this section we focus on the greedy algorithms and we use those to select a subset of visual features. We implemented the greedy algorithms and the construction of the matrices required in the functions (9) in c++, using eigen for the computation of the log-determinant and the smallest eigenvalue. For numerical reasons, rather than computing the determinant and taking the logarithm, we directly compute the log-determinant from the Cholesky decomposition of the matrix. For the computation of the smallest eigenvalue we use eigen's svd function.

Our feature selection approach is used as an add-on to a visual-inertial pipeline similar to the one described in [68]. Our VIN pipeline estimates the navigation state (robot pose, velocity, and IMU biases) using the structureless visual model and the pre-integrated IMU model described in [68]. The entire implementation is based on the GTSAM optimization library [80]. Our implementation differs from [68] in three important ways. First, in this paper we use the iSAM2 algorithm within a fixed-lag smoothing approach; we marginalize out states outside a smoothing horizon of 6s, which helps bounding latency and memory requirements. Second, we do not adopt SVO as visual front-end; in this simulations we do not need a front-end as we simulate landmark observations, while in the following section we describe a simple real-world front-end. Finally, rather than feeding to the VIN estimator all available measurements, we use the feature selection algorithms described in this paper to select a small set of informative visual observations.

In this section we evaluate two main aspects of our approach. First, we show that a clever selection of the features does actually impact VIN accuracy. Second, we show that the lazy evaluation approach discussed in Section IV-C speeds up the computation of the greedy solution. We use two metrics for accuracy: the *absolute translation error*, which is the Euclidean distance between the estimated position and the actual position, and the *relative translation error*, which computes the Euclidean norm of the difference between the estimated translation between time k and time $k + 1$ and the actual translation. Indeed the relative translation error quantifies how quickly the estimate drifts from the ground truth. Since absolute positions are not observable in visual-inertial

odometry, the relative error is a more reliable performance metric. When useful, we report absolute and relative rotation errors (defined in analogy with the translation ones).

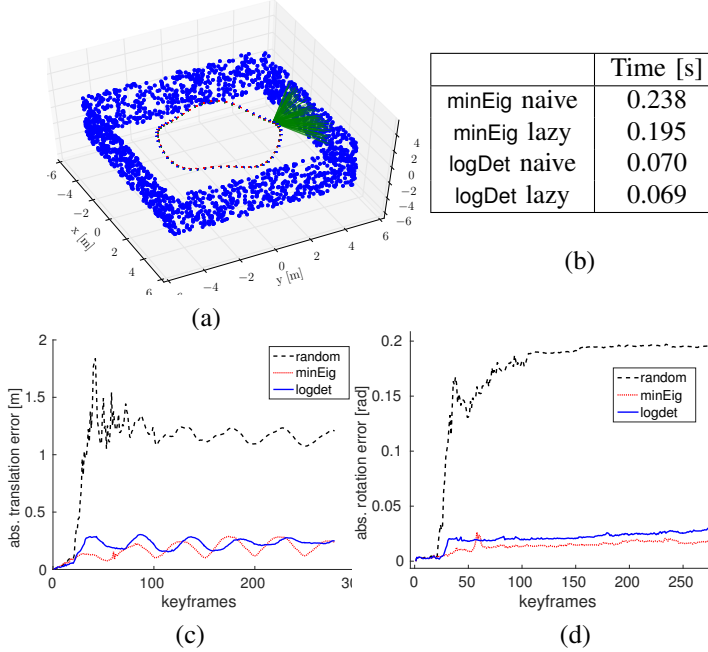


Fig. 2. Simulation results: (a) simulated environment, (b) table with CPU times for different implementations of the greedy algorithms, (c) absolute translation errors, (d) absolute rotation errors.

Results. We simulate 50 Monte Carlo runs; in each run we add random noise to the acceleration, gyroscope, and camera measurements. To make the simulation realistic, the statistics about measurement noise are identical to the ones used in the real tests of Section V-C. In each run, the robot performs VIN and, at each camera frames, it selects $\kappa = 20$ visual features out of all the features in the field of view. We compare three feature selection strategies. The greedy selection resulting from Algorithm 1 with the eigenvalue objective f_λ (label: minEig), the greedy selection with the log-determinant cost f_{\det} (label: logDet), and a random selection which randomly draws κ of the available features (label: random).

Fig. 2(c)-(d) show the absolute translation and absolute rotation errors, averaged over 50 Monte Carlo runs. From the figure is clear that a clever selection of the features, resulting from logDet or minEig, deeply impacts performance in VIN. Our techniques largely improve estimation errors, when compared against the random selection; both approaches result in similar performance. From the figure we note that the absolute errors have some oscillations: this is a consequence of the fact that the trajectory is circular; in general, this stresses the fact that absolute metrics may be poor indicators of performance in visual-inertial odometry. In this case, the relative error metrics confirm the results of Fig. 2(c)-(d): the average translation and rotation errors are given in Table I; in parenthesis we report the error reduction percentage with respect to the random baseline.

Fig. 2(b) reports the CPU time required for feature selection. The figure considers both cost functions (logDet and minEig) and

Technique	Rel. Translation Error [m]	Rel. Rotation Error [rad]
random	0.0103	0.0049
minEig	0.0064 (-37%)	0.0025 (-48%)
logDet	0.0053 (-48%)	0.0018 (-63%)

TABLE I
RELATIVE TRANSLATION AND ROTATION ERRORS FOR THE SIMULATED TESTS OF SECTION V-B (AVERAGE OVER 50 MONTE CARLO RUNS)

compares timing when using our lazy evaluation, as described in Algorithm 1, against a naive implementation of the greedy algorithm that always tests the marginal gain of every feature (i.e., for which the stopping condition in line 12 of Algorithm 1 is disabled). The naive greedy (without lazy evaluation) always results in κN objective evaluations. When using lazy evaluation, the number of objective evaluation depends on the tightness of the upper bounds used in Algorithm 1. From Fig. 2(b), we see that the advantage of using the lazy evaluation is marginal for the log-determinant cost; this is not surprising, since the Hadamard's inequality of Proposition 3 usually gives a fairly loose bound. On the other hand, the advantage of using the lazy evaluation is significant for the minEig, resulting in a reduction of the computational time of 20%. The average CPU time required by Algorithm 1 (with lazy evaluation) to select $\kappa = 20$ features is 0.069s for logDet and 0.195s for minEig. While these timing may be already acceptable for applications, there are large margins to speed up computation: we postpone these considerations to Section VI.

C. Real Tests: Agile Navigation on Micro Aerial Vehicles

In this section we show that our feature selection approach enhances VIN performance in real-world navigation problems with micro-aerial vehicles (MAVs).

Testing setup. We use the EuRoC benchmark [11] for our evaluation. The EuRoC datasets are collected with an AscTech Firefly hex-rotor helicopter equipped with a VI (stereo) visual-inertial sensor. The camera records stereo images at resolution 752×480 and framerate 20Hz; IMU data is collected at 200Hz. We refer to [11] for a technical description of the datasets. In this context we only remark that the datasets contain test instances at increasing levels of complexity, collected in a machine hall environment and in a smaller Vicon room. In our tests, the measurement variances, as well as the intrinsic and extrinsic calibration parameters match exactly the one specified in the dataset. The most relevant parameters used in our tests are given in Table II; in the front-end we used openCV's *goodFeaturesToTrack* for feature detection and the Lucas-Kanade method for feature tracking; as input to the detector we specify a minimum quality level for the features and a desired number of features to extract (N). From these N features our selector has to retain $\kappa = 10$ features that will be used by the back-end. In this sense, feature detection and tracking at the front-end are pre-attentive mechanisms: they work on a large set of features, which are later filtered out by our feature selector. The feature selector uses a predictive horizon of 3s; in practice, the future pose estimates along the horizon can be computed from the control inputs, by integrating the dynamics of the vehicle (Section III-B). Since the control inputs are not available in the EuRoC dataset, we compute the future poses by attaching ground truth motion

increments to the current pose estimate. The only assumption in doing so is that the control loop and the estimation quality are good enough to track a desired set of future poses; this is the case in VIN in which the short-term drift is small.

	Parameter name	Value
Front-end	Nr. features to detect (N)	100
	Minimum quality level	0.001
	Time between keyframes	0.2s
Back-end	Smoothing window	6s
	iSAM2 iterations	1
Feature selector	Nr. features to select (κ)	10
	Horizon	3s

TABLE II
VIN AND FEATURE SELECTION PARAMETERS

Techniques. We compare four VIN approaches. The first two VIN approaches use the minEig and the logDet selectors proposed in this paper. The third approach uses a selector that picks the κ features with highest quality (i.e., highest score in *goodFeaturesToTrack*). This selector is commonly used in VIN and only accounts for the appearance of the visual features; we denote it with the label “quality”, following openCV’s terminology. The fourth technique is a VIN approach using 200 features (selected as the ones with largest score in *goodFeaturesToTrack*) and is used to have a reference performance for the case in which the VIN system has less stringent computational constraints (label: no-selection).

In order to compute the tracking probabilities p_l , we modified openCV’s *goodFeaturesToTrack* in order to have access to the features’ scores. Then, we mapped the scores to probabilities in $[0, 1]$, such that more distinguishable features have higher tracking probabilities p_l .

Results: Accuracy. Fig. 3 shows the performance of the compared techniques on all the 11 EuRoC datasets. The EuRoC benchmark includes datasets of different levels of complexity, with the difficult datasets being challenging for standard VIN pipeline due to the fast motion of the MAV. In this section we show that we can obtain accurate position estimation with as few as $\kappa = 10$ features; this budget is enforced for each frame; for instance, if we are tracking r features from the previous frame, then in the current frame we can only retain $\kappa - r$ features.

Fig. 3(a) compares the VIN performance using the relative translation errors as metric. The figure confirms that the difficult datasets tend to have larger translation errors. Moreover, it shows that the proposed techniques, minEig and logDet, lead to smallest errors compared to the baseline quality. Clearly, the technique no-selection, which uses 20x more features, leads to the smaller errors. To better appreciate the advantage of minEig and logDet with respect to quality, Fig. 3(b) shows the relative improvement, i.e., the relative translation error reduction, of the two techniques with respect to quality. The figure shows that the proposed feature selectors result in much smaller drift across all but one datasets. The average error reduction is larger than 20% and overcomes 40% in the datasets MH_02_easy, MH_05_difficult, and V2_01_easy. In particular, in the dataset MH_05_difficult the estimate resulting from the quality-based feature selection diverged after a sharp turn, while our techniques were able to ensure

accurate pose estimation. The dataset V1_03_difficult is the only one in which the proposed techniques have slightly worse performance. We noticed that in datasets with severe motion blur the advantage of the proposed techniques may vary, and this is due to the fact that we are using a simplistic model for the tracking probabilities p_l . For completeness, Fig. 3(c) reports the absolute translation error as a percentage of the trajectory traveled; this is another common metric for VIN. We notice that no-selection has excellent performance, while using 200 features (average error accumulation is 0.17% of the trajectory length). The approach no-selection exhibits similar or smaller estimation errors with respect to related techniques benchmarked on the EuRoC datasets in [81].¹⁰ Moreover, the proposed techniques, logDet and minEig, are able to ensure an average error accumulation of 0.42% and 0.46%, respectively, while using only 10 features!

To get a better intuition behind the large performance boost induced by the proposed techniques, we report few snapshots produced by our pipeline in Fig. 4. Each sub-figure shows, for the current frame, the tracked features (green squares with the optical flow vector), the available features (red crosses), and the features selected (yellow circles) by (a) quality, (b) logDet, and (c) minEig. The frames are captured during a sharp left turn from the MH_03_medium dataset. The quality selector simply picks the most distinguishable features, resulting in many features selected on the right-hand side of the image; these features are of scarce utility: they will soon disappear from the field of view due to the motion of the MAV. On the other hand, logDet and minEig are predictive and they leverage the knowledge of the immediate motion of the platform; therefore they tend to discard features that fall outside the field of view and select features on the left-hand side of the image.

Results: Timing and Trade-offs. Fig. 5 reports the average CPU time required by the VIN back-end for all techniques and datasets. For the proposed techniques, the back-end time includes both the CPU time spent on feature selection and the time spent for estimation in VIN (factor graph optimization with iSAM2). The figure shows that logDet is able to reduce the back-end time by 30–40% in most scenarios, with respect to no-selection; in particular, the average time for no-selection is 57ms while the average time for logDet is 35ms. The CPU time of quality is even smaller, at the cost of degraded performance (Fig. 3). Consistent with the Monte Carlo analysis, in our current implementation logDet is faster than minEig, which implies a back-end time larger than 100ms in our tests.

Fig. 6 provides a more detailed breakdown of the back-end time for increasing number of selected features and for each of the EuRoC datasets. The figure focuses on the logDet approach, which currently has been already shown to be remarkably faster than the minEig approach. The blue portion of each bar reports the time spent on feature selection, while the red portion corresponds to the time spent on factor graph optimization

¹⁰The interested reader can compare Fig. 2(a) in [81] with Fig. 3(c) in the present paper, noting that we report errors as percentage of the distance traveled and 10 out of 11 of the EuRoC datasets are less than 100m long [11]; furthermore, we notice that [81] reported systematic failures on 2 of the datasets, while we are able to successfully complete all the datasets.

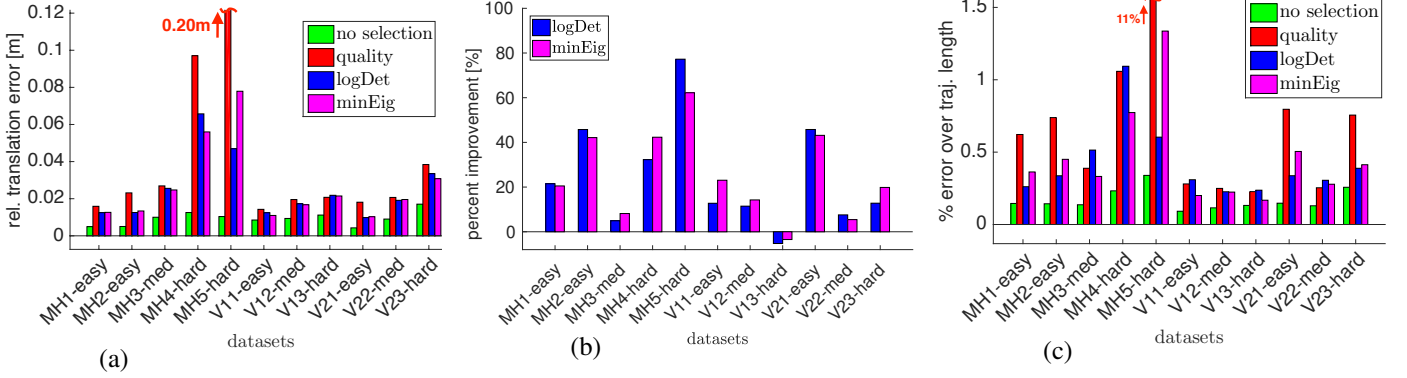


Fig. 3. Accuracy for the compared techniques on the 11 EuRoC MAV datasets. (a) Relative translation error; (b) Relative improvement (relative translation error reduction) of the proposed techniques with respect to the quality baseline; (c) Translation error as percentage of the overall trajectory length.

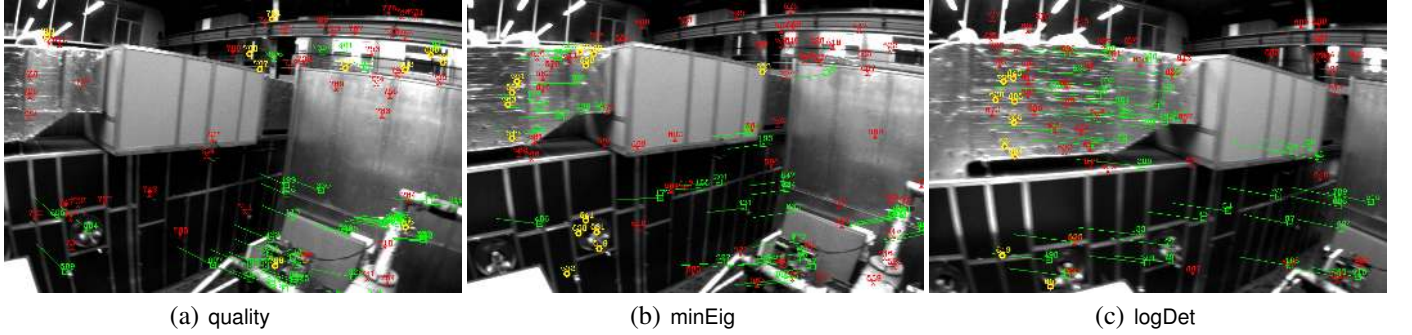


Fig. 4. Snapshots of the feature selection performed by the techniques quality, minEig, and logDet during a sharp left turn. Features tracked from previous frames are shown as green squares (with the corresponding optical flow vectors), the newly detected features are shown as red crosses, and the selected features are shown as yellow circles. We note that quality only selects the features from their appearance, and chooses many features on the right-hand side of the frames: these features will soon fall out of the field of view due to the sharp turn.

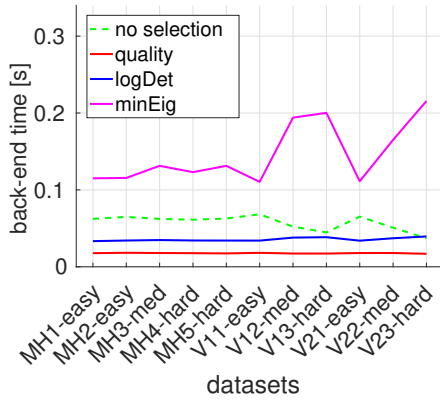


Fig. 5. CPU time for the back-end (including feature selection) for the compared techniques on the EuRoC datasets.

in iSAM2; the sum of these two times corresponds to the overall back-end time reported in Fig. 5. The figure compares the timing results for the proposed approach versus no-selection for increasing number of selected features (in all cases, no-selection retains all the tracked features). The figure provides few useful insights. First, the computational advantage in using the proposed attention mechanism is more noticeable when the number of selected features is small. In particular, when selecting 5-10 features we typically have a large computational saving by using the proposed approach. However, when the number of features approaches the overall number of available

measurements (100 in our tests), the computational advantage can be negligible. This is due to the fact that our current implementation of the feature selection is relatively slow and indeed its computational cost is often comparable with the cost of running iSAM2. In Section VI we discuss extensions that have the potential of making the selection time negligible. The second insight is that the computational advantage of the proposed approach is more evident in easy datasets, while in some of the difficult datasets (e.g., V12_med, V13_hard, V22_med, V23_hard) the computational advantage becomes marginal or inexistent. This mismatch results from the fact that in the easy datasets the front-end is able to track many features (typically all the features that we set as upper bound) which implies that the no-selection approach requires more time to perform estimation. In the more difficult scenarios, instead, the front-end typically tracks less features hence resulting in faster estimation (and degraded performance). Note that while in the difficult scenarios the computational advantage of our approach may be limited, the accuracy boost resulting from our approach still suggests its use over quality.

Fig. 7 provides the computation/performance trade-off for increasing number of features and for each of the EuRoC datasets, using the logDet selector. The back-end time, shown as a dashed red line, corresponds to the sum of the feature selection time and the iSAM2 estimation time. As expected the time increases with the number of selected features. The

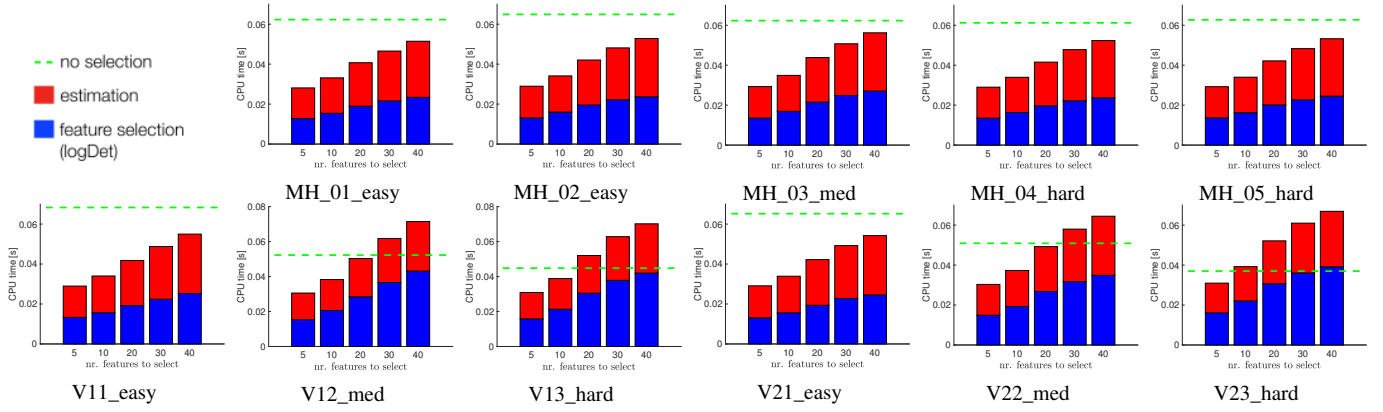


Fig. 6. CPU time breakdown for the logDet selector for increasing number of selected features and for each of the EuRoC datasets. The blue portion of each bar reports the time spent on feature selection (Algorithm 1), while the red portion corresponds to the time spent on factor graph optimization in iSAM2; the sum of these two times corresponds to the overall back-end time. The back-end time is compared against the CPU time required by no-selection, shown as a dashed green line.

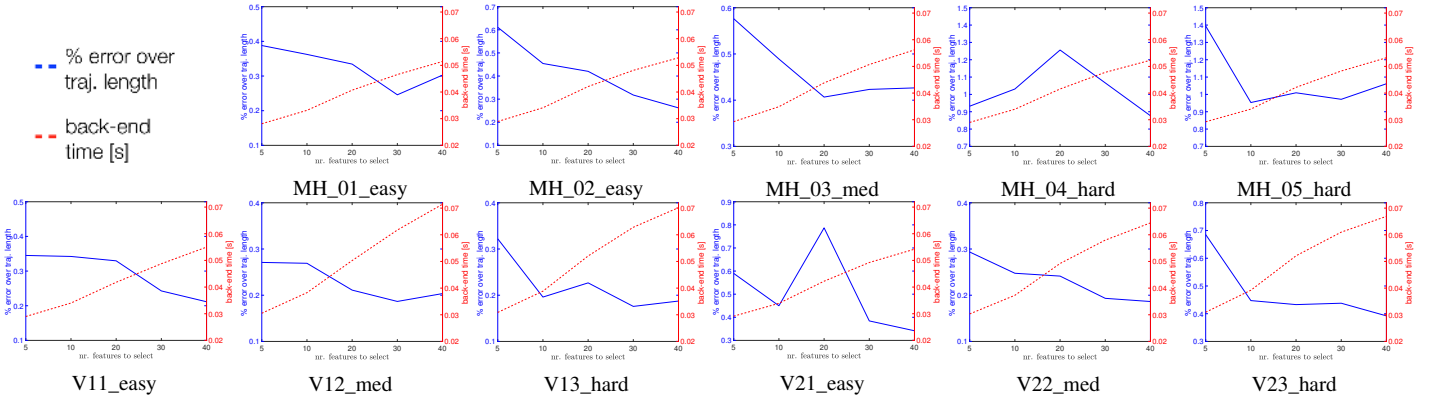


Fig. 7. Performance/computation trade-off for the logDet selector for increasing number of selected features and for each of the EuRoC datasets. The back-end time, shown as a dashed red line, corresponds to the sum of the feature selection time and the iSAM2 estimation time. The figure also reports the absolute translation error as a percentage of the trajectory traveled, shown as a solid blue line. Results are averaged over 5 runs performed on each dataset.

figure also reports the absolute translation error as a percentage of the trajectory traveled, shown as a solid blue line. The error typically decreases when increasing the number of selected features. Results are averaged over 5 runs performed on each dataset. The trade-off plots in Fig. 7 can be used to decide the number of features to use to attain a desired level of accuracy or given an upper bound on the time that can be spent performing estimation at the back-end.

VI. CONCLUSION AND FUTURE WORK

This work provides an attention mechanism for visual-inertial navigation. This mechanism takes the form of a feature selector, which retains the most informative visual features detected by the VIN front-end (pre-attentive process) and feeds them to the estimation back-end. We proposed two algorithms for feature selection. Both algorithms enjoy four desirable qualities: they are predictive in nature, in that they are aware of the motion of the robot in the immediate future; they are task-driven, since they select a set of features that minimize the VIN estimation error; they are greedy, hence efficient and easy to implement; they come with performance guarantees that bound their sub-optimality. We demonstrated our feature selection extensively on both realistic Monte Carlo simulations

and real-world data collected by a micro aerial vehicle. The experiments suggest that the feature selection seriously impacts VIN performance; the use of the proposed techniques reduces the estimation error in easy datasets, and enables accurate estimation in difficult datasets in which standard approaches would fail on a limited budget of visual features. This work opens many avenues for future investigation.

Computational improvements. The first avenue for future work consists in reducing the computational time of feature selection. Two main ideas can make the feature selection time negligible. The first stems from the observation that the greedy algorithm is trivially parallelizable: the marginal gain of each feature can be computed independently; leveraging this fact alone would result in large computational savings. The second idea is to use sparse matrix manipulation to compute the determinant and the smallest eigenvalues; our current implementation uses dense matrices. Another interesting avenue consists in including a learning mechanism to improve feature selection. A learning-based method may allow capturing more complex sensory-motor dynamics and may improve the feature selection over time, potentially adjusting to changing environments and time-varying robot dynamics.

Task-driven perception. A second avenue for future work consists in extending our attention framework. We plan to explore two paths. First, while (1) minimizes the localization uncertainty subject to a feature budget, one may also consider a “dual” problem in which one minimizes the number of features to be used, while satisfying a desired localization performance. From the technical standpoint, this alternative formulation can be tackled in a similar manner and in both cases greedy algorithms have sub-optimality guarantees. This alternative formulation would provide a grounded answer to the question: *how much visual information is needed to navigate at a desired accuracy?* The second avenue consists in extending our attention framework to other tasks: for instance, *how many visual features does the robot need to sense in order to avoid crashing into nearby obstacles?* We believe these are necessary steps towards a task-driven perception theory, that can enable autonomy on robots with tight resource constraints.

ACKNOWLEDGMENTS

The authors gratefully acknowledge Christian Forster for the useful discussions and for providing the VIN simulator, and the anonymous reviewers for providing invaluable comments to improve the quality of this manuscript.

APPENDIX

A. Linear Imu Model

In this appendix we provide explicit expressions for the matrices and vectors appearing in the IMU model (14).

Given the velocity \mathbf{v}_k and the position \mathbf{t}_k of the robot at time k , we can get \mathbf{v}_{k+1} and \mathbf{v}_{k+1} by simple Euler integration using the acceleration \mathbf{a}_k :

$$\mathbf{v}_{k+1} = \mathbf{v}_k + \mathbf{a}_k \delta \quad (38)$$

$$\mathbf{t}_{k+1} = \mathbf{t}_k + \mathbf{v}_k \delta + \frac{1}{2} \mathbf{a}_k \delta^2 \quad (39)$$

By induction, the velocity and position at time $j > k$ are:

$$\begin{aligned} \mathbf{v}_j &= \mathbf{v}_k + \sum_{i=k}^{j-1} \mathbf{a}_i \delta \\ \mathbf{t}_j &= \mathbf{t}_k + \sum_{i=k}^{j-1} \mathbf{v}_i \delta + \frac{1}{2} \sum_{i=k}^{j-1} \mathbf{a}_i \delta^2 \\ &\quad (\text{substituting } \mathbf{v}_i) \\ &= \mathbf{t}_k + \sum_{i=k}^{j-1} (\mathbf{v}_k + \sum_{h=k}^{i-1} \mathbf{a}_h \delta) \delta + \frac{1}{2} \sum_{i=k}^{j-1} \mathbf{a}_i \delta^2 \\ &\quad (\text{moving } \mathbf{v}_k \text{ outside the sum}) \\ &= \mathbf{t}_k + (j-k) \mathbf{v}_k \delta + \sum_{i=k}^{j-1} \sum_{h=k}^{i-1} \mathbf{a}_h \delta^2 \\ &\quad + \frac{1}{2} \sum_{i=k}^{j-1} \mathbf{a}_i \delta^2 \\ &\quad (\text{developing the double sum for } i = k, \dots, j-1 \\ &\quad \text{and noting that the summand } i = k \text{ vanishes}) \\ &= \mathbf{t}_k + (j-k) \mathbf{v}_k \delta \\ &\quad + \underbrace{\mathbf{a}_k \delta^2}_{i=k+1} + \underbrace{(\mathbf{a}_k + \mathbf{a}_{k+1}) \delta^2}_{i=k+2} \\ &\quad + \underbrace{(\mathbf{a}_k + \mathbf{a}_{k+1} + \mathbf{a}_{k+2}) \delta^2}_{i=k+3} + \dots \\ &\quad + \underbrace{(\mathbf{a}_k + \mathbf{a}_{k+1} + \dots + \mathbf{a}_{j-2}) \delta^2}_{i=j-1} \end{aligned} \quad (40)$$

$$+ \underbrace{(\mathbf{a}_k + \mathbf{a}_{k+1} + \mathbf{a}_{k+2}) \delta^2}_{i=k+3} + \dots \quad (41)$$

$$+ \underbrace{(\mathbf{a}_k + \mathbf{a}_{k+1} + \dots + \mathbf{a}_{j-2}) \delta^2}_{i=j-1} \quad (42)$$

$$+ \frac{1}{2} \sum_{i=k}^{j-1} \mathbf{a}_i \delta^2 \quad (43)$$

Now we note that in the expanded sum in lines (40)-(42) the term \mathbf{a}_k appears $j-k-1$ times, the term \mathbf{a}_{k+1} appears $j-k-2$ times, and so on, till the term \mathbf{a}_{j-2} which appears just once. Therefore, we can write lines (40)-(42) succinctly as a single sum, leading to:

$$\begin{aligned} \mathbf{t}_j &= \mathbf{t}_k + (j-k) \mathbf{v}_k \delta + \sum_{i=k}^{j-2} (j-i-1) \mathbf{a}_i \delta^2 + \frac{1}{2} \sum_{i=k}^{j-1} \mathbf{a}_i \delta^2 \\ &\quad (\text{putting last two terms together}) \\ &= \mathbf{t}_k + (j-k) \mathbf{v}_k \delta + \frac{1}{2} \mathbf{a}_{j-1} \delta^2 + \sum_{i=k}^{j-2} (j-i-\frac{1}{2}) \mathbf{a}_i \delta^2 \\ &\quad (\text{simplifying}) \\ &= \mathbf{t}_k + (j-k) \mathbf{v}_k \delta + \sum_{i=k}^{j-1} (j-i-\frac{1}{2}) \mathbf{a}_i \delta^2 \end{aligned}$$

Defining $\delta_{kj} \doteq (j-k) \delta$ and substituting \mathbf{a}_k from (10):

$$\begin{aligned} \mathbf{v}_j &= \mathbf{v}_k + \sum_{i=k}^{j-1} (\mathbf{R}_i (\tilde{\mathbf{a}}_i - \mathbf{b}_k - \boldsymbol{\eta}_i) + \mathbf{g}) \delta \\ &= \mathbf{v}_k + \mathbf{g} \delta_{kj} - (\sum_{i=k}^{j-1} \mathbf{R}_i \delta) \mathbf{b}_k + \sum_{i=k}^{j-1} \mathbf{R}_i \tilde{\mathbf{a}}_i \delta \\ &\quad - \sum_{i=k}^{j-1} \mathbf{R}_i \boldsymbol{\eta}_i \delta \\ \mathbf{t}_j &= \mathbf{t}_k + \mathbf{v}_k \delta_{kj} + \sum_{i=k}^{j-1} (j-i-\frac{1}{2}) (\mathbf{R}_i (\tilde{\mathbf{a}}_i - \mathbf{b}_k - \boldsymbol{\eta}_i) + \mathbf{g}) \delta^2 \\ &= \mathbf{t}_k + \mathbf{v}_k \delta_{kj} + \sum_{i=k}^{j-1} (j-i-\frac{1}{2}) \mathbf{g} \delta^2 \\ &\quad - (\sum_{i=k}^{j-1} (j-i-\frac{1}{2}) \mathbf{R}_i \delta^2) \mathbf{b}_k + \sum_{i=k}^{j-1} (j-i-\frac{1}{2}) \mathbf{R}_i (\tilde{\mathbf{a}}_i - \boldsymbol{\eta}_i) \delta^2 \end{aligned} \quad (44)$$

Let us now define the following vectors:

$$\begin{aligned} \mathbf{z}_{kj}^{\mathbf{v}} &\doteq \mathbf{g} \delta_{kj} + \sum_{i=k}^{j-1} \mathbf{R}_i \tilde{\mathbf{a}}_i \delta \\ \boldsymbol{\eta}_{kj}^{\mathbf{v}} &\doteq \sum_{i=k}^{j-1} \mathbf{R}_i \boldsymbol{\eta}_i \delta \\ \mathbf{z}_{kj}^{\mathbf{t}} &\doteq \sum_{i=k}^{j-1} (j-i-\frac{1}{2}) \mathbf{g} \delta^2 \\ &\quad + \sum_{i=k}^{j-1} (j-i-\frac{1}{2}) \mathbf{R}_i \tilde{\mathbf{a}}_i \delta^2 \\ \boldsymbol{\eta}_{kj}^{\mathbf{t}} &\doteq + \sum_{i=k}^{j-1} (j-i-\frac{1}{2}) \mathbf{R}_i \boldsymbol{\eta}_i \delta^2 \end{aligned} \quad (45)$$

Using this notation we rewrite eq. (44) (putting position first) and adding the random walk random model on the bias:

$$\begin{aligned} \mathbf{t}_j &= \mathbf{t}_k + \mathbf{v}_k \delta_{kj} - \left(\sum_{i=k}^{j-1} (j-i-\frac{1}{2}) \mathbf{R}_i \delta^2 \right) \mathbf{b}_k + \mathbf{z}_{kj}^{\mathbf{t}} - \boldsymbol{\eta}_{kj}^{\mathbf{t}} \\ \mathbf{v}_j &= \mathbf{v}_k - (\sum_{i=k}^{j-1} \mathbf{R}_i \delta) \mathbf{b}_k + \mathbf{z}_{kj}^{\mathbf{v}} - \boldsymbol{\eta}_{kj}^{\mathbf{v}} \\ \mathbf{b}_j &= \mathbf{b}_k - \boldsymbol{\eta}_{kj}^{\mathbf{b}} \end{aligned} \quad (46)$$

In order to write (46) in compact matrix form, we define:

$$\mathbf{N}_{kj} \doteq \sum_{i=k}^{j-1} (j-i-\frac{1}{2}) \mathbf{R}_i \delta^2 \quad (47)$$

$$\mathbf{M}_{kj} \doteq \sum_{i=k}^{j-1} \mathbf{R}_i \delta \quad (48)$$

which allows rewriting (46) succinctly as:

$$\begin{aligned} \mathbf{z}_{kj}^{\mathbf{t}} &= \mathbf{t}_j - \mathbf{t}_k - \mathbf{v}_k \delta_{kj} + \mathbf{N}_{kj} \mathbf{b}_k + \boldsymbol{\eta}_{kj}^{\mathbf{t}} \\ \mathbf{z}_{kj}^{\mathbf{v}} &= \mathbf{v}_j - \mathbf{v}_k + \mathbf{M}_{kj} \mathbf{b}_k + \boldsymbol{\eta}_{kj}^{\mathbf{v}} \\ \mathbf{z}_{kj}^{\mathbf{b}} &= \mathbf{b}_j - \mathbf{b}_k + \boldsymbol{\eta}_{kj}^{\mathbf{b}} \end{aligned} \quad (49)$$

where $\mathbf{z}_{kj}^{\mathbf{b}} = \mathbf{0}_3$ is the expected change in the bias.

Let us now define the following matrices and vectors:

$$\begin{aligned} \mathbf{A}_{kj} &= \begin{bmatrix} \mathbf{0}_{9 \times 9} & \dots & -\mathbf{I}_3 & -\mathbf{I}_3 \delta_{kj} & \mathbf{N}_{kj} \\ \mathbf{0} & & \mathbf{0} & -\mathbf{I}_3 & \mathbf{M}_{kj} \\ \mathbf{0} & & \mathbf{0} & \mathbf{0} & -\mathbf{I}_3 \end{bmatrix} \begin{bmatrix} \mathbf{I}_9 & \mathbf{0}_{9 \times 9} & \dots \end{bmatrix} \\ \mathbf{z}_{kj}^{\text{IMU}} &= \begin{bmatrix} \mathbf{z}_{kj}^{\mathbf{t}} \\ \mathbf{z}_{kj}^{\mathbf{v}} \\ \mathbf{z}_{kj}^{\mathbf{b}} \end{bmatrix} \quad \boldsymbol{\eta}_{kj}^{\text{IMU}} = \begin{bmatrix} \boldsymbol{\eta}_{kj}^{\mathbf{t}} \\ \boldsymbol{\eta}_{kj}^{\mathbf{v}} \\ \boldsymbol{\eta}_{kj}^{\mathbf{b}} \end{bmatrix} \end{aligned} \quad (50)$$

Using (50), we finally rewrite our model (49) as:

$$\mathbf{z}_{kj}^{\text{IMU}} = \mathbf{A}_{kj} \mathbf{x}_{k:k+H} + \boldsymbol{\eta}_{kj}^{\text{IMU}} \quad (51)$$

To fully characterize the linear measurement model (51) we only have to compute the covariance of the noise $\boldsymbol{\eta}_{kj}^{\text{IMU}}$, which is given by:

$$\text{cov}(\boldsymbol{\eta}_{kj}^{\text{IMU}}) = \begin{bmatrix} \sigma_{\text{IMU}}^2 \mathbf{C} \mathbf{C}^\top & \mathbf{0}_{6 \times 3} \\ \mathbf{0}_{3 \times 6} & \text{cov}(\boldsymbol{\eta}_{kj}^{\text{b}}) \end{bmatrix} \quad (52)$$

where \mathbf{C} includes the coefficient matrices of the noise in (46):

$$\mathbf{C} = \begin{bmatrix} (j-k-\frac{1}{2})\mathbf{R}_k \delta^2 & (j-k-\frac{3}{2})\mathbf{R}_{k+1} \delta^2 & \dots & \frac{1}{2}\mathbf{R}_{j-1} \delta^2 \\ \mathbf{R}_k \delta & \mathbf{R}_{k+1} \delta & \dots & \mathbf{R}_{j-1} \delta \end{bmatrix}$$

Using the fact that any rotation matrix satisfies $\mathbf{R}_k^\top \mathbf{R}_k = \mathbf{I}_3$, the matrix $\mathbf{C} \mathbf{C}^\top$ can be computed simply as:

$$\mathbf{C} \mathbf{C}^\top = \begin{bmatrix} \left(\sum_{i=k}^{j-1} (j-i-\frac{1}{2})^2 \right) \delta^4 \mathbf{I}_3 & \left(\sum_{i=k}^{j-1} (j-i-\frac{1}{2}) \right) \delta^3 \mathbf{I}_3 \\ \left(\sum_{i=k}^{j-1} (j-i-\frac{1}{2}) \right) \delta^3 \mathbf{I}_3 & (j-k-1) \delta^2 \mathbf{I}_3 \end{bmatrix}.$$

B. Proof of Proposition 1

The information matrix of the joint state (22) is additive in the measurements, hence the information matrices of the joint states $[\mathbf{x}_{k:k+H} \ \mathbf{p}_{l_1}]$ and $[\mathbf{x}_{k:k+H} \ \mathbf{p}_{l_2}]$ given the predicted visual measurements to landmarks l_1 and l_2 can be respectively written as:

$$\boldsymbol{\Omega}_{k:H}^{(l_1)} = \sum_{\tau=k}^{k_1} \boldsymbol{\Omega}_{\tau}^{(l_1)} \quad \boldsymbol{\Omega}_{k:H}^{(l_2)} = \sum_{\tau=k}^{k_2} \boldsymbol{\Omega}_{\tau}^{(l_2)} \quad (53)$$

where $\boldsymbol{\Omega}_{\tau}^{(l_1)}$ (resp. $\boldsymbol{\Omega}_{\tau}^{(l_2)}$) is the contribution to the information matrix corresponding to the measurement of landmark 1 (resp. 2) at time τ . Since the proposition assumes that the future measurements are identical, it follows $\sum_{\tau=k}^{k_1} \boldsymbol{\Omega}_{\tau}^{(l_1)} = \sum_{\tau=k}^{k_1} \boldsymbol{\Omega}_{\tau}^{(l_2)}$; this, combined with the fact that $k_2 > k_1$ implies:

$$\boldsymbol{\Omega}_{k:H}^{(l_2)} = \boldsymbol{\Omega}_{k:H}^{(l_1)} + \sum_{\tau=k_1+1}^{k_2} \boldsymbol{\Omega}_{\tau}^{(l_2)} \succeq \boldsymbol{\Omega}_{k:H}^{(l_1)} \quad (54)$$

Now we only have to prove that the Schur complement preserves the ordering $\boldsymbol{\Omega}_{k:H}^{(l_2)} \succeq \boldsymbol{\Omega}_{k:H}^{(l_1)}$, since Δ_{l_1} and Δ_{l_2} are simply the Schur complements of $\boldsymbol{\Omega}_{k:H}^{(l_1)}$ and $\boldsymbol{\Omega}_{k:H}^{(l_2)}$, respectively. For this purpose, we first observe that

$$\boldsymbol{\Omega}_{k:H}^{(l_2)} \succeq \boldsymbol{\Omega}_{k:H}^{(l_1)} \implies \left(\boldsymbol{\Omega}_{k:H}^{(l_2)} \right)^{-1} \preceq \left(\boldsymbol{\Omega}_{k:H}^{(l_1)} \right)^{-1} \quad (55)$$

Moreover, we make explicit the block structure of the two matrices as follows:

$$\boldsymbol{\Omega}_{k:H}^{(l_1)} \doteq \begin{bmatrix} \mathbf{A}_1 & \mathbf{B}_1 \\ \mathbf{B}_1^\top & \mathbf{C}_1 \end{bmatrix} \quad \boldsymbol{\Omega}_{k:H}^{(l_2)} \doteq \begin{bmatrix} \mathbf{A}_2 & \mathbf{B}_2 \\ \mathbf{B}_2^\top & \mathbf{C}_2 \end{bmatrix} \quad (56)$$

where the upper-left blocks (\mathbf{A}_1 and \mathbf{A}_2) correspond to entries of the information matrix associated to the states $\mathbf{x}_{k:k+H}$ and the bottom-right blocks (\mathbf{C}_1 and \mathbf{C}_2) correspond to the landmark states we want to marginalize.

Now we note that using standard block inversion for the block matrix $\boldsymbol{\Omega}_{k:H}^{(l_1)}$ we obtain:

$$\left(\boldsymbol{\Omega}_{k:H}^{(l_1)} \right)^{-1} = \begin{bmatrix} (\mathbf{A}_1 - \mathbf{B}_1 \mathbf{C}_1^{-1} \mathbf{B}_1^\top)^{-1} & \star \\ \star & \star \end{bmatrix} \quad (57)$$

where we denoted with “ \star ” blocks which are irrelevant for the following derivation. Combining the inequality (55) with the block inverse (57) we get:

$$\begin{bmatrix} (\mathbf{A}_2 - \mathbf{B}_2 \mathbf{C}_2^{-1} \mathbf{B}_2^\top)^{-1} & \star \\ \star & \star \end{bmatrix} \preceq \begin{bmatrix} (\mathbf{A}_1 - \mathbf{B}_1 \mathbf{C}_1^{-1} \mathbf{B}_1^\top)^{-1} & \star \\ \star & \star \end{bmatrix} \quad (58)$$

Since diagonal blocks of positive semidefinite matrices are also semidefinite, eq. (58) implies $(\mathbf{A}_2 - \mathbf{B}_2 \mathbf{C}_2^{-1} \mathbf{B}_2^\top)^{-1} \preceq (\mathbf{A}_1 - \mathbf{B}_1 \mathbf{C}_1^{-1} \mathbf{B}_1^\top)^{-1}$ hence:

$$(\mathbf{A}_2 - \mathbf{B}_2 \mathbf{C}_2^{-1} \mathbf{B}_2^\top) \succeq (\mathbf{A}_1 - \mathbf{B}_1 \mathbf{C}_1^{-1} \mathbf{B}_1^\top). \quad (59)$$

Comparing the block structure in (56) with the description in eqs. (22)-(23), we realize that $(\mathbf{A}_1 - \mathbf{B}_1 \mathbf{C}_1^{-1} \mathbf{B}_1^\top) = \Delta_{l_1}$ and $(\mathbf{A}_2 - \mathbf{B}_2 \mathbf{C}_2^{-1} \mathbf{B}_2^\top) = \Delta_{l_2}$ hence (59) implies $\Delta_{l_1} \preceq \Delta_{l_2}$ concluding the proof.

C. Proof of Corollary 5

The proof relies on the inequality (31) for i chosen to be the smallest eigenvalue. From the Weyl inequality [76], it follows $\lambda_j(\mathbf{M} + \Delta) \geq \lambda_{\min}(\mathbf{M})$, for all j . Using this fact, it follows that the minimum in (31) is attained by $\lambda_{\min}(\mathbf{M} + \Delta)$. Therefore, the inequality (31) becomes:

$$|\lambda_{\min}(\mathbf{M}) - \lambda_{\min}(\mathbf{M} + \Delta)| \leq \|\Delta \mathbf{v}_{\min}\| \quad (60)$$

From the positive definiteness of \mathbf{M} and Δ (which implies $\lambda_{\min}(\mathbf{M}) \geq 0$ and $\lambda_{\min}(\mathbf{M} + \Delta) \geq 0$), and from the Weyl inequality, it follows $|\lambda_{\min}(\mathbf{M}) - \lambda_{\min}(\mathbf{M} + \Delta)| = \lambda_{\min}(\mathbf{M} + \Delta) - \lambda_{\min}(\mathbf{M})$, which substituted in (60) leads to (32).

D. Proof of Proposition 12

In order to show that the submodularity ratio (36) does not vanish, we show that its numerator is bounded away from zero. To do so, we consider a single summand in (36):

$$(a) \doteq f(\mathbf{L} \cup \{e\}) - f(\mathbf{L}) = \lambda_{\min}(\boldsymbol{\Omega}_{\mathbf{L}} + \Delta_e) - \lambda_{\min}(\boldsymbol{\Omega}_{\mathbf{L}}) \quad (61)$$

where $\boldsymbol{\Omega}_{\mathbf{L}} \doteq \bar{\boldsymbol{\Omega}}_{k:k+H} + \sum_{l \in \mathbf{L}} \Delta_l$. Our task is to prove that (61) is different from zero. To do so, we substitute the eigenvalue with its definition through the Rayleigh quotient:

$$\begin{aligned} (a) &= \min_{\|\boldsymbol{\mu}\|=1} \boldsymbol{\mu}^\top (\boldsymbol{\Omega}_{\mathbf{L}} + \Delta_e) \boldsymbol{\mu} - \min_{\|\boldsymbol{\nu}\|=1} \boldsymbol{\nu}^\top (\boldsymbol{\Omega}_{\mathbf{L}}) \boldsymbol{\nu} \\ &\quad (\text{calling } \bar{\boldsymbol{\mu}} \text{ the minimizer of the first summand}) \\ &= \bar{\boldsymbol{\mu}}^\top (\boldsymbol{\Omega}_{\mathbf{L}} + \Delta_e) \bar{\boldsymbol{\mu}} - \min_{\|\boldsymbol{\nu}\|=1} \boldsymbol{\nu}^\top (\boldsymbol{\Omega}_{\mathbf{L}}) \boldsymbol{\nu} \\ &\quad (\text{since } \boldsymbol{\nu} = \bar{\boldsymbol{\mu}} \text{ is suboptimal for the second summand}) \\ &\geq \bar{\boldsymbol{\mu}}^\top (\boldsymbol{\Omega}_{\mathbf{L}} + \Delta_e) \bar{\boldsymbol{\mu}} - \bar{\boldsymbol{\mu}}^\top (\boldsymbol{\Omega}_{\mathbf{L}}) \bar{\boldsymbol{\mu}} \\ &\quad (\text{simplifying and substituting the expression of } \Delta_e \text{ from (23)}) \\ &= \bar{\boldsymbol{\mu}}^\top \Delta_e \bar{\boldsymbol{\mu}} = \bar{\boldsymbol{\mu}}^\top \mathbf{F}_e^\top (\mathbf{I} - \mathbf{E}_e (\mathbf{E}_e^\top \mathbf{E}_e)^{-1} \mathbf{E}_e^\top) \mathbf{F}_e \bar{\boldsymbol{\mu}} \\ &\quad (\text{defining the idempotent matrix } \mathbf{Q}_e \doteq (\mathbf{I} - \mathbf{E}_e (\mathbf{E}_e^\top \mathbf{E}_e)^{-1} \mathbf{E}_e^\top)) \\ &= \bar{\boldsymbol{\mu}}^\top \mathbf{F}_e^\top \mathbf{Q}_e \mathbf{F}_e \bar{\boldsymbol{\mu}} = \bar{\boldsymbol{\mu}}^\top \mathbf{F}_e^\top \mathbf{Q}_e \mathbf{Q}_e \mathbf{F}_e \bar{\boldsymbol{\mu}} = \|\mathbf{Q}_e \mathbf{F}_e \bar{\boldsymbol{\mu}}\|^2 \end{aligned}$$

Now we write \mathbf{E}_e in terms of its 3×3 blocks:

$$\mathbf{E}_e = \begin{bmatrix} \mathbf{E}_{e0} \\ \mathbf{E}_{e1} \\ \vdots \\ \mathbf{E}_{eH} \end{bmatrix} \quad (62)$$

Moreover, we recall that F_e has the following block structure:

$$F_e = \left[\begin{array}{ccc|ccc|ccc|c} -E_{e0} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ 0 & 0 & 0 & -E_{e1} & 0 & 0 & 0 & 0 & 0 & \cdots \\ 0 & 0 & 0 & 0 & 0 & 0 & -E_{e2} & 0 & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{array} \right]$$

where we noted that the nonzero blocks in F_e are the same (up-to-sign) as the ones in E_e (c.f. the coefficient matrices in (19)). It follows that $F_e \bar{\mu}$ can be written explicitly as:

$$F_e \bar{\mu} = \begin{bmatrix} -E_{e0} \bar{\mu}_0 \\ -E_{e1} \bar{\mu}_1 \\ \vdots \\ -E_{eH} \bar{\mu}_H \end{bmatrix} \quad (63)$$

Now we observe that Q_e is an orthogonal projector onto the null space of E_e , and the null space of Q_e is spanned by the columns of E_e . Therefore, any vector v that falls in the null space of Q_e can be written as a linear combination of the columns of E_e :

$$Q_e v = 0 \Leftrightarrow v = E_e w \quad (64)$$

with $w \in \mathbb{R}^3$. By comparison with (62), we note that $F_e \bar{\mu}$ can be written as $E_e w$ if and only if $\bar{\mu}_1 = \bar{\mu}_2 = \dots = \bar{\mu}_H$. Therefore, if $\bar{\mu}_i \neq \bar{\mu}_j$ for some i, j , then the vector $F_e \bar{\mu}$ cannot be in the null space of Q_e , and the lower bound (62) must be greater than zero, concluding the proof.

REFERENCES

- [1] M. Potter, B. Wyble, C. Hagmann, and E. McCourt, "Detecting meaning in RSVP at 13 ms per picture," *Attention, Perception, & Psychophysics*, vol. 76, no. 2, pp. 270–279, 2014.
- [2] M. Carrasco, "Visual attention: The past 25 years," *Vision Research*, vol. 51, pp. 1484–1525, 2011.
- [3] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [4] S. Pillai, S. Ramalingam, and J. Leonard, "High-performance and tunable stereo reconstruction," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2016.
- [5] R. Mur-Artal, J. Montiel, and J. Tardós, "ORB-SLAM: A versatile and accurate monocular SLAM system," *IEEE Trans. Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [6] NVIDIA GeForce Website, "Geforce GTX TITAN X specifications." [Online]. Available: <http://www.geforce.com/hardware/desktop-gpus/geforce-gtx-titan-x/specifications>
- [7] Wikipedia, "List of cpu power dissipation figures." [Online]. Available: https://en.wikipedia.org/wiki/List_of_CPU_power_dissipation_figures
- [8] Z. Zhang, A. Suleiman, L. Carlone, V. Sze, and S. Karaman, "Visual-inertial odometry on chip: An algorithm-and-hardware co-design approach," in *Robotics: Science and Systems (RSS)*, 2017, (pdf) (web), highlighted in the MIT News: (web).
- [9] P. Cavanagh, "Visual cognition," *Vision Research*, vol. 51, no. 13, pp. 1538 – 1551, 2011, vision Research 50th Anniversary Issue: Part 2. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0042698911000381>
- [10] L. Carlone and S. Karaman, "Attention and anticipation in fast visual-inertial navigation," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2017, pp. 3886–3893, extended arxiv preprint: 1610.03344 (pdf).
- [11] M. Burri, J. Nikolic, P. Gohl, T. Schneider, J. Rehder, S. Omari, M. Achtelik, and R. Siegwart, "The EuRoC micro aerial vehicle datasets," *Intl. J. of Robotics Research*, 2016.
- [12] A. Borji and L. Itti, "State-of-the-art in visual attention modeling," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 35, no. 1, 2013.
- [13] B. Scholl, "Objects and attention: the state of the art," *Cognition*, vol. 80, no. 1, pp. 1–46, 2001.
- [14] D. Caduff and S. Timpf, "On the assessment of landmark salience for human navigation," *Cogn. Process*, vol. 9, pp. 249–267, 2008.
- [15] J. Moran and J. Desimone, "Selective attention gates visual processing in the extrastriate cortex," *Science*, vol. 229, no. 4715, pp. 782–784, 1985.
- [16] J. Wolfe, "Guided search 2.0 - a revised model of visual search," *Psychon Bull. Rev.*, vol. 1, no. 2, pp. 202–238, 1994.
- [17] H. Spekreijse, "Pre-attentive and attentive mechanisms in vision. perceptual organization and dysfunction," *Vision Research*, vol. 40, no. 10-12, pp. 1179–1182, 2000.
- [18] R. Sim and G. Dudek, "Learning and evaluating visual features for pose estimation," in *Intl. Conf. on Computer Vision (ICCV)*, 1999, pp. 1217–1222.
- [19] V. Peretroukhin, L. Clement, M. Giamou, and J. Kelly, "PROBE: Predictive robust estimation for visual-inertial navigation," in *IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2015.
- [20] N. Ouerhani, A. Bur, and H. Hügli, "Visual attention-based robot self-localization," in *ECMR*, 2005, pp. 8–13.
- [21] P. Newman and K. Ho, "Slam-loop closing with visually salient features," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2005, pp. 635–642.
- [22] P. Sala, R. Sim, A. Shokoufandeh, and S. Dickinson, "Landmark selection for vision-based navigation," *IEEE Trans. Robotics*, vol. 22, no. 2, pp. 334–349, 2006.
- [23] C. Siagian and L. Itti, "Biologically-inspired robotics vision monte-carlo localization in the outdoor environment," in *IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2007, pp. 1723–1730.
- [24] S. Frintrop and P. Jensfelt, "Attentional landmarks and active gaze control for visual SLAM," *IEEE Trans. Robotics*, vol. 24, no. 5, pp. 1054–1065, 2008.
- [25] S. Hochdorfer and C. Schlegel, "Landmark rating and selection according to localization coverage: Addressing the challenge of lifelong operation of slam in service robots," in *IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2009, pp. 382–387.
- [26] H. Strasdat, C. Stachniss, and W. Burgard, "Which landmark is useful? learning selection policies for navigation in unknown environments," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2009, pp. 1410–1415.
- [27] M. Chli and A. Davison, "Active matching for visual tracking," *Robotics and Autonomous Systems*, vol. 57, no. 12, pp. 1173–1187, Dec. 2009.
- [28] A. Handa, M. Chli, H. Strasdat, and A. J. Davison, "Scalable active matching," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2010.
- [29] J. Jang, K. Won, and S. Jung, "Geometric feature selection for vehicle pose estimation on dynamic road scenes," in *5th Intl. Conf. on Ubiquitous and Future Networks*, 2013.
- [30] Z. Shi, Z. Liu, X. Wu, and W. Xu, "Feature selection for reliable data association in visual SLAM," *Machine Vision and Applications*, vol. 24, pp. 667–682, 2013.
- [31] A. Oliva and A. Torralba, "Modeling the shape of the scene: a holistic representation of the spatial envelope," *Intl. J. of Computer Vision*, vol. 42, pp. 145–175, 2001.
- [32] A. Torralba, A. Oliva, M. Castelhan, and J. Henderson, "Contextual guidance of attention in natural scenes: the role of global features on object search," *Psychological Review*, vol. 113, no. 4, pp. 766–786, 2006.
- [33] V. Mnih, N. Heess, A. Graves, and K. Kavukcuoglu, "Recurrent models of visual attention," in *Advances in Neural Information Processing Systems (NIPS)*, 2014, pp. 2204–2212.
- [34] K. Xu, J. L. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *ArXiv preprint: 1502.03044*, 2016.
- [35] I. Cvišić and I. Petrović, "Stereo odometry based on careful feature selection and tracking," in *Proc. of the European Conference on Mobile Robots (ECMR)*, 2015.
- [36] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, "Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age," *IEEE Trans. Robotics*, vol. 32, no. 6, pp. 1309–1332, 2016, arxiv preprint: 1606.05830, (pdf).
- [37] A. Davison, "Active search for real-time vision," in *Intl. Conf. on Computer Vision (ICCV)*, Oct 2005.

- [38] R. Lerner, E. Rivlin, and I. Shimshoni, "Landmark selection for task-oriented navigation," *IEEE Trans. Robotics*, vol. 23, no. 3, pp. 494–505, 2007.
- [39] B. Mu, A. Agha-mohammadi, L. Paull, M. Graham, J. How, and J. Leonard, "Two-stage focused inference for resource-constrained collision-free navigation," in *Robotics: Science and Systems (RSS)*, 2015.
- [40] K. Wu, T. Do, L. C. Carrillo-Arce, and S. I. Roumeliotis, "On the VINS resource-allocation problem for a dual-camera, small-size quadrotor," in *Intl. Sym. on Experimental Robotics (ISER)*, 2016, pp. 538–549.
- [41] G. Zhang and P. Vela, "Good features to track for visual slam," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [42] D. Kottas, R. DuToit, A. Ahmed, C. Guo, G. Georgiou, R. Li, and S. Roumeliotis, "A resource-aware vision-aided inertial navigation system for wearable and portable computers," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2014, pp. 6336–6343.
- [43] S. Joshi and S. Boyd, "Sensor selection via convex optimization," *IEEE Trans. Signal Processing*, vol. 57, pp. 451–462, 2009.
- [44] C. Giraud and B. Jouvencel, "Sensor selection: A geometrical approach," in *IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, vol. 2, 1995, pp. 1410–1415.
- [45] A. Krause, A. Singh, and C. Guestrin, "Near-optimal sensor placements in gaussian processes: Theory, efficient algorithms and empirical studies," *Journal of Machine Learning Research*, vol. 9, pp. 235–284, 2008.
- [46] F. Bian, D. Kempe, and R. Govindan, "Utility based sensor selection," in *5th Int. Conf. Information Processing Sensor Networks*, 2006, pp. 11–18.
- [47] M. Shamaiah, S. Banerjee, and H. Vikalo, "Greedy sensor selection: leveraging submodularity," in *IEEE Conf. on Decision and Control (CDC)*, 2010.
- [48] M. Vitus, W. Zhang, A. Abate, J. Hu, and C. Tomlin, "On efficient sensor scheduling for linear dynamical systems," *Automatica*, vol. 48, pp. 2482–2493, 2012.
- [49] H. Zhang, R. Ayoub, and S. Sundaram, "Sensor selection for optimal filtering of linear dynamical systems: Complexity and approximation," in *IEEE Conf. on Decision and Control (CDC)*, 2015.
- [50] S. Jawaid and S. Smith, "Submodularity and greedy algorithms in sensor scheduling for linear dynamical systems," *Automatica*, vol. 61, pp. 282–288, 2015.
- [51] V. Tzoumas, A. Jadbabaie, and G. Pappas, "Sensor placement for optimal kalman filtering: Fundamental limits, submodularity, and algorithms," in *American Control Conference*, 2016.
- [52] T. Summers, F. Cortesi, and J. Lygeros, "On submodularity and controllability in complex dynamical networks," *IEEE Transactions on Control of Network Systems*, vol. 3, no. 1, pp. 91–101, 2016.
- [53] A. Mourikis and S. Roumeliotis, "A multi-state constraint Kalman filter for vision-aided inertial navigation," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, April 2007, pp. 3565–3572.
- [54] A. Davison, I. Reid, N. Molton, and O. Stasse, "MonoSLAM: Real-time single camera SLAM," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 29, no. 6, pp. 1052–1067, Jun 2007.
- [55] M. Bloesch, S. Omari, M. Hutter, and R. Siegwart, "Robust visual inertial odometry using a direct EKF-based approach," in *IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2015.
- [56] E. Jones and S. Soatto, "Visual-inertial navigation, mapping and localization: A scalable real-time causal approach," *Intl. J. of Robotics Research*, vol. 30, no. 4, Apr 2011.
- [57] J. Hesch, D. Kottas, S. Bowman, and S. Roumeliotis, "Camera-imu-based localization: Observability analysis and consistency improvement," *Intl. J. of Robotics Research*, vol. 33, no. 1, pp. 182–201, 2014.
- [58] A. Mourikis and S. Roumeliotis, "A dual-layer estimator architecture for long-term localization," in *Proc. of the Workshop on Visual Localization for Mobile Platforms at CVPR*, Anchorage, Alaska, June 2008.
- [59] G. Sibley, L. Matthies, and G. Sukhatme, "Sliding window filter with application to planetary landing," *J. of Field Robotics*, vol. 27, no. 5, pp. 587–608, 2010.
- [60] T.-C. Dong-Si and A. Mourikis, "Motion tracking with fixed-lag smoothing: Algorithm consistency and analysis," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2011.
- [61] S. Leutenegger, S. Lynen, M. Bosse, R. Siegwart, and P. Furgale, "Keyframe-based visual-inertial slam using nonlinear optimization," *Intl. J. of Robotics Research*, 2015.
- [62] M. Bryson, M. Johnson-Roberson, and S. Sukkarieh, "Airborne smoothing and mapping using vision and inertial sensors," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2009, pp. 3143–3148.
- [63] V. Indelman, S. Williams, M. Kaess, and F. Dellaert, "Information fusion in navigation systems via factor graph based incremental smoothing," *Robotics and Autonomous Systems*, vol. 61, no. 8, pp. 721–738, August 2013.
- [64] S. Shen, *Autonomous Navigation in Complex Indoor and Outdoor Environments with Micro Aerial Vehicles*. PhD Thesis, University of Pennsylvania, 2014.
- [65] N. Keivan, A. Patron-Perez, and G. Sibley, "Asynchronous adaptive conditioning for visual-inertial SLAM," in *Intl. Sym. on Experimental Robotics (ISER)*, 2014.
- [66] A. Patron-Perez, S. Lovegrove, and G. Sibley, "A spline-based trajectory representation for sensor fusion and rolling shutter cameras," *Intl. J. of Computer Vision*, February 2015.
- [67] T. Lupton and S. Sukkarieh, "Visual-inertial-aided navigation for high-dynamic motion in built environments without initial conditions," *IEEE Trans. Robotics*, vol. 28, no. 1, pp. 61–76, Feb 2012.
- [68] C. Forster, L. Carlone, F. Dellaert, and D. Scaramuzza, "On-manifold preintegration theory for fast and accurate visual-inertial navigation," *IEEE Trans. Robotics*, vol. 33, no. 1, pp. 1–21, 2016, arxiv preprint: 1512.02363, (pdf), technical report GT-IRIM-CP&R-2015-001.
- [69] A. Das and D. Kempe, "Submodular meets spectral: Greedy algorithms for subset selection, sparse approximation and dictionary selection," in *Intl. Conf. on Machine Learning (ICML)*, 2011, pp. 1057–1064.
- [70] J. Farrell, *Aided Navigation: GPS with High Rate Sensors*. McGraw-Hill, 2008.
- [71] R. Tron, L. Carlone, F. Dellaert, and K. Daniilidis, "Rigid components identification and rigidity enforcement in bearing-only localization using the graph cycle basis," in *American Control Conference*, 2015, (pdf).
- [72] L. Carlone, Z. Kira, C. Beall, V. Indelman, and F. Dellaert, "Eliminating conditionally independent sets in factor graphs: A unifying perspective based on smart factors," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2014.
- [73] L. Carlone, P. F. Alcantarilla, H. Chiu, K. Zsolt, and F. Dellaert, "Mining structure fragments for smart bundle adjustment," in *British Machine Vision Conf. (BMVC)*, 2014, accepted as oral presentation (acceptance rate 7.7%), (pdf) (supplemental material: (pdf)).
- [74] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge University Press, 2004.
- [75] R. Horn and C. Johnson, *Matrix Analysis*. Cambridge University Press, 1985.
- [76] I. C. F. Ipsen and B. Nadler, "Refined perturbation bounds for eigenvalues of hermitian and non-hermitian matrices," *SIAM J. Matrix Analysis*, vol. 31, no. 1, 2009.
- [77] G. Nemhauser, L. Wolsey, and M. Fisher, "An analysis of approximations for maximizing submodular set functions," *Mathematical Programming*, vol. 14, no. 1, pp. 265–294, 1978.
- [78] A. Das, A. Dasgupta, and R. Kumar, "Selecting diverse features via spectral regularization," in *Advances in Neural Information Processing Systems (NIPS)*, vol. 2, 2012, pp. 1583–1591.
- [79] D. Rosen, L. Carlone, A. Bandeira, and J. Leonard, "SE-Sync: A certifiably correct algorithm for synchronization over the Special Euclidean group," in *Intl. Workshop on the Algorithmic Foundations of Robotics (WAFR)*, 2016, extended arxiv preprint: 1611.00128, (pdf) (pdf) (code).
- [80] F. Dellaert, "Factor graphs and GTSAM: A hands-on introduction," Georgia Institute of Technology, Tech. Rep. GT-RIM-CP&R-2012-002, September 2012.
- [81] K. Sun, K. Mohta, B. Pfrommer, M. Watterson, S. Liu, Y. Mulgaonkar, C. Taylor, and V. Kumar, "Robust stereo visual inertial odometry for fast autonomous flight," *IEEE Robotics and Automation Letters*, vol. 3, no. 2, pp. 965–972, 2018.