

UNIVERSITÉ CLERMONT AUVERGNE
École Universitaire de Physique et d'Ingénierie

AprilTag-based Trajectory Tracking for the LIMOS Robot

Second-year Master's project

Automation, Robotics track: Artificial Perception and Robotics

Authors:

Joao Pedro MARTINS DO LAGO REIS
Yann Kelvem DA SILVA RAMOS

Supervisor:

Dr. Sébastien LENGAGNE

Defense Date:

March 02, 2026

Aubière, Clermont-Ferrand

Title

Subtitle

Author

Abstract

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Contents

1	Introduction	2
1.1	Motivation	2
1.1.1	testesubsub	3
1.2	Objectives	3
1.2.1	General objective	3
1.2.2	Specific objectives	3
2	Fundamentacao teorica	4
2.1	Câmera RGB	4
2.2	AprilTag	5
2.2.1	Decodificação confiável: códigos, rotações e distância de Hamming	6
2.2.2	Detecção robusta: limiarização adaptativa e extração do quadrilátero	7
2.2.3	Por que a pose é estimável: geometria projetiva e homografia	7
2.2.4	AprilTag 3 e confiabilidade sob generalização de layouts	8
3	Otmizacao e trajetoria	9
3.0.1	Subsubsection	10
4	Metodologia	11
5	Results	12
6	Discussion	13
7	Conclusion	14
	References	15

1 Introduction

Being able to follow a path once and repeat it optimally is a practical skill for mobile robots, especially in indoor environments where GPS is not available and tasks are repetitive by nature. Thus, this is the central idea behind Teach and Replay[1], the robot first performs a guided route while it registers the markers looking for a positioning reference of each tag and then uses this record of positions to reproduce the same route optimally.

In this practical exercise, this concept was implemented on a LIMO mobile robot using its integrated RGB camera and AprilTag [2] fiducial markers positioned along a closed circuit with a start and end. During the teach phase, the robot completes a first turn while detecting the markers and saves the information from the positions of the robot in relation to the tags to represent the executed movement. After pressing the Y key of the Xbox controller or another controller's triangle, the system switches to the replay phase, where the robot tries to follow the learned trajectory. The goal is simple: complete the circuit, reducing the path deviation and maintaining a safe distance from the tags, but a restriction is imposed on the user of at least always seeing two tags in the camera.

In order to develop the solution in a controlled way and evaluate its behavior in real conditions, a simulation workflow was made for reality. Therefore, it started in the Gazebo Ignition, where the perception and control components have been tested and validated. We then transfer the same workflow to the real LIMO robot, where external noise such as lighting variation, motion blur, and wheel slippage naturally increase the challenge of detection quality for tracking performance. This report therefore documents not only the final system but also the changes that occurred when moving from simulation to a real platform.

The rest of the report presents the requirements of the problem, justifying that it will be useful a tracking system trajectory, and also details the methodology proposed for the method Teach and Replay, both for the part of the validation in simulation and for the implementation in the real case, defines the evaluation metrics and discusses the results obtained with a direct comparison between simulated and real executions.

1.1 Motivation

The motivation to do the work was that it will serve to help students from related areas of programming and robotics at Polytechnique de Clermont-Ferrand to use the system in a practical work, so they can understand how the simple camera perception part, using markers, can be applied directly to mobile robotics, understand the complexity of the transformations of positions used and how, with a simple camera and markers, it is possible to trace a trajectory and, with optimization, arrive at the best possible trajectory.

1.1.1 testesubsub

The motivation to do the work was that it will serve to help students from related areas of programming and robotics at Polytechnique de Clermont-Ferrand to use the system in a practical work, so they can understand how the simple camera perception part, using markers, can be applied directly to mobile robotics, understand the complexity of the transformations of positions used and how, with a simple camera and markers, it is possible to trace a trajectory and, with optimization, arrive at the best possible trajectory.

1.2 Objectives

1.2.1 General objective

1.2.2 Specific objectives

2 Fundamentacao teorica

2.1 Câmera RGB

A câmera RGB embarcada no robô LIMO constitui o principal sensor exteroceptivo utilizado neste projeto, fornecendo uma sequência de imagens 2D a partir das quais são extraídas observações geométricas (cantos de marcadores, contornos e, conseqüentemente, pose relativa). Do ponto de vista de modelagem, assume-se que o sistema óptico pode ser aproximado por um modelo de câmera pinhole[3] com distorção radial/tangencial, uma hipótese amplamente empregada em visão computacional quando o objetivo é relacionar pontos 3D no espaço a coordenadas 2D na imagem com precisão suficiente para estimação de pose.

Conforme ilustrado na Figura 2.1, um ponto 3D no referencial da câmera projeta-se no plano de imagem segundo o modelo pinhole.

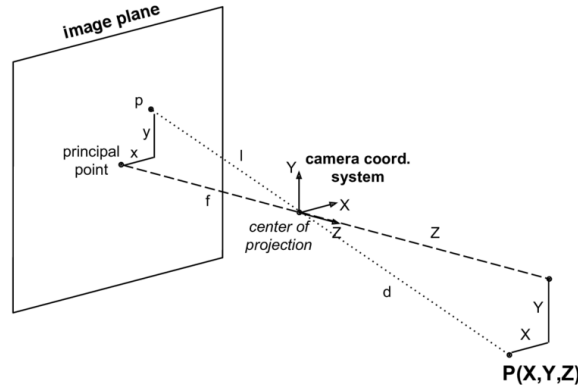


Figure 2.1: Modelo pinhole: projeção de um ponto 3D $P(X, Y, Z)$ no plano de imagem, indicando o centro de projeção e o ponto principal.

Seja um ponto no referencial da câmera $\mathbf{p}_c = [X_c \ Y_c \ Z_c]^\top$, com $Z_c > 0$. Sua projeção ideal[3] (sem distorção) no plano normalizado é dada por

$$x_n = \frac{X_c}{Z_c}, \quad y_n = \frac{Y_c}{Z_c}. \quad (2.1)$$

A passagem para coordenadas de pixel é determinada pela matriz intrínseca K ,

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_n \\ y_n \\ 1 \end{bmatrix}, \quad (2.2)$$

onde f_x, f_y são as distâncias focais em unidades de pixel e (c_x, c_y) é o centro principal. Na prática, lentes reais introduzem distorção, e a relação entre coordenadas normalizadas ideais (x_n, y_n) e coordenadas distorcidas (x_d, y_d) é frequentemente modelada por termos radiais e tangenciais. Definindo $r^2 = x_n^2 + y_n^2$, um modelo comum é

$$x_d = x_n(1 + k_1 r^2 + k_2 r^4 + k_3 r^6) + 2p_1 x_n y_n + p_2(r^2 + 2x_n^2), \quad (2.3)$$

$$y_d = y_n(1 + k_1 r^2 + k_2 r^4 + k_3 r^6) + p_1(r^2 + 2y_n^2) + 2p_2 x_n y_n, \quad (2.4)$$

e então (u, v) são obtidos substituindo (x_d, y_d) em (2.2). A relevância dessas equações para o presente trabalho decorre do fato de que a estimação de pose baseada em marcadores planares (como AprilTags) depende da precisão na localização de cantos e, portanto, da fidelidade do mapeamento entre o espaço projetivo e o plano de imagem. Mesmo pequenas distorções não modeladas podem introduzir vies sistemático nos cantos detectados, o que se propaga para erros em rotação e translação.

Além da geometria, propriedades físicas do sensor influenciam diretamente a qualidade da percepção. Em particular, ruído de intensidade, compressão e variações de exposição alteram o contraste local utilizado em binarização/adaptação de limiar, enquanto o desfoque de movimento (motion blur) reduz a nitidez de bordas e prejudica a extração de contornos. Um modelo simplificado do motion blur pode ser escrito como uma convolução espacial

$$I_{\text{blur}} = I * h, \quad (2.5)$$

onde I é a imagem ideal e h é um kernel de borrimento dependente da velocidade relativa e do tempo de exposição. No contexto de detecção de marcadores, essa degradação afeta principalmente a estimação subpixel dos cantos, reduzindo a repetibilidade das poses obtidas entre frames.

Finalmente, para tarefas de navegação, a câmera deve ser interpretada como um sensor que produz medições com incerteza. Assim, as observações de cantos (e, por consequência, a pose estimada) podem ser vistas como variáveis aleatórias. Um modelo usual assume erro gaussiano aditivo em pixels,

$$\hat{\mathbf{u}}_i = \mathbf{u}_i + \boldsymbol{\epsilon}_i, \quad \boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}, \Sigma_u), \quad (2.6)$$

o que justifica, posteriormente, o uso de refinamento por minimização de erro de reprojeção (na estimação de pose) e a adoção de estratégias de controle robustas no seguimento de trajetória.

Em síntese, o modelo pinhole com intrínsecos e distorção ((2.1)–(2.4)), aliado às considerações de ruído e borrimento ((2.5)–(2.6)), fornece a base teórica para compreender como a câmera RGB do LIMO influencia a detecção de marcadores e a qualidade das estimativas geométricas utilizadas no restante do sistema.

2.2 AprilTag

AprilTags são marcadores fiduciais planares projetados para fornecer, a partir de uma única imagem, duas informações fundamentais para robótica móvel: (i) uma identificação discreta (ID) e (ii) uma estimativa geométrica de pose 6-DoF do marcador em relação à câmera. A proposta do sistema não é maximizar capacidade de dados, mas sim garantir detecção estável e baixa taxa de falsos positivos sob variações de escala, rotação e iluminação, o que as torna adequadas como *landmarks* artificiais em tarefas de localização e navegação [2]. O detector foi redesenhado no AprilTag 2 com foco explícito em eficiência e robustez computacional [4], e o AprilTag 3 generaliza o uso de layouts e discute critérios de complexidade para preservar confiabilidade mesmo ao ampliar o espaço de padrões possíveis [5].

Para contextualizar a estrutura de uma tag e o papel da borda e das células internas na decodificação, recomenda-se inserir na Figura 2.2 um exemplo de AprilTag (família e ID) destacando a borda e a região de dados. Em seguida, para conectar a teoria ao pipeline perceptivo, recomenda-se inserir na Figura 2.3 um fluxograma curto do detector (binarização adaptativa, extração de contornos/quadriláteros, retificação e decodificação), conforme discutido em [4], [5].

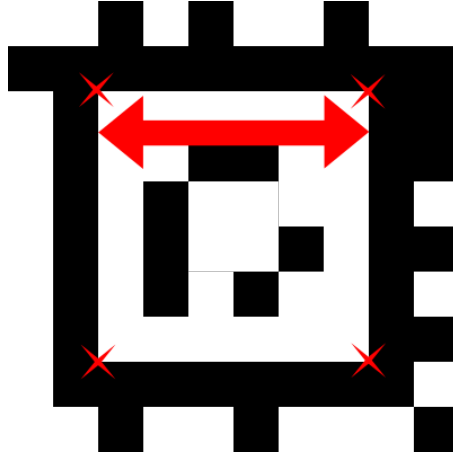


Figure 2.2: Estrutura de uma AprilTag: borda de alto contraste e região de dados (payload) usada na identificação.

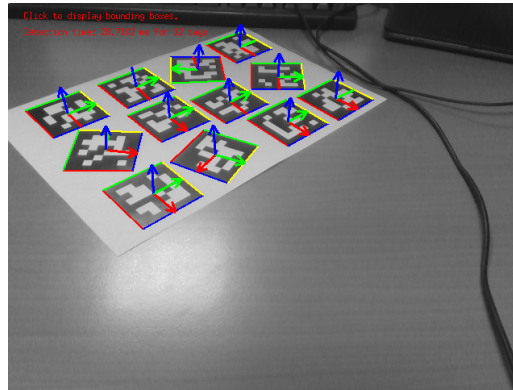


Figure 2.3: Pipeline conceitual de detecção: imagem → limiarização adaptativa → extração de candidatos (quads) → retificação por homografia → leitura de bits e validação do ID.

2.2.1 Decodificação confiável: códigos, rotações e distância de Hamming

A confiabilidade do ID decorre do projeto de famílias de códigos que se mantêm separáveis mesmo quando a tag é observada sob rotações de 0° , 90° , 180° e 270° . Considerando palavras binárias $a, b \in \{0, 1\}^n$, a separação entre dois códigos é medida pela distância de Hamming,

$$d_H(a, b) = \sum_{k=1}^n \mathbb{I}\{a_k \neq b_k\}, \quad (2.7)$$

e a geração da família impõe uma distância mínima (incluindo as rotações como equivalências), o que reduz ambiguidades de orientação e diminui a probabilidade de padrões naturais serem aceitos como IDs válidos [2], [4]. Para uma distância mínima d , a capacidade teórica clássica de correção por decisão de vizinho mais próximo é

$$t = \left\lfloor \frac{d-1}{2} \right\rfloor, \quad (2.8)$$

embora, em aplicações robóticas, uma prática frequente seja rejeitar leituras com erros para diminuir ainda mais falsos positivos [4]. Essa combinação (família bem separada e validação estrita) é a base para tratar AprilTags como marcadores “fiáveis” em ambientes reais.

Figure 2.4: Exemplo de detecção em imagem: quadrilátero estimado, ID decodificado e (opcionalmente) eixos de pose sobrepostos.

2.2.2 Detecção robusta: limiarização adaptativa e extração do quadrilátero

O detector precisa localizar geometricamente a tag antes de decodificar bits. Em cenários reais, um limiar global falha com sombras e iluminação não uniforme; por isso, o AprilTag 2 usa limiarização adaptativa baseada em estatísticas locais (mínimo e máximo) para obter uma binarização que preserve transições preto-branco relevantes [4]. Uma forma canônica de expressar esse princípio é definir, para uma vizinhança $\Omega(x, y)$, um limiar

$$T(x, y) = \frac{\max_{(i,j) \in \Omega(x,y)} I(i, j) + \min_{(i,j) \in \Omega(x,y)} I(i, j)}{2}, \quad I_b(x, y) = \mathbb{I}\{I(x, y) \geq T(x, y)\}, \quad (2.9)$$

onde I é a imagem em tons de cinza e I_b a imagem binária. Em seguida, etapas eficientes de segmentação (por exemplo, por componentes conexas) e agrupamento de bordas permitem construir candidatos a quadriláteros (quads) que satisfaçam restrições geométricas compatíveis com a projeção de um quadrado planar [4]. No AprilTag 3, mantém-se a estrutura do pipeline e discute-se ainda o controle de desempenho por decimação e ajustes voltados a aumentar *recall* em tags pequenas [5].

Para ilustrar a saída perceptiva em nível de imagem (o que o detector realmente “vê”), recomenda-se inserir na Figura 2.4 um frame real ou do simulador com o contorno do quad e o ID sobreposto; se disponível, incluir também os eixos de pose desenhados sobre a tag (isso ajuda a conectar diretamente com a seção de metodologia/controlre).

2.2.3 Por que a pose é estimável: geometria projetiva e homografia

Uma AprilTag é um alvo planar de geometria conhecida; logo, ao observar seus quatro cantos na imagem, a relação entre o plano da tag e o plano de imagem é modelada por uma transformação projetiva (homografia). Assumindo câmera pinhole com matriz intrínseca K , um ponto $\mathbf{X} = [X \ Y \ 0]^\top$ no plano da tag projeta-se como

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K \begin{bmatrix} \text{camera} R_{\text{tag}} & \text{camera} t_{\text{tag}} \end{bmatrix} \begin{bmatrix} X \\ Y \\ 0 \\ 1 \end{bmatrix} = K \begin{bmatrix} r_1 & r_2 & t \end{bmatrix} \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix}, \quad (2.10)$$

o que define a homografia

$$H = K \begin{bmatrix} r_1 & r_2 & t \end{bmatrix}, \quad (2.11)$$

com r_1, r_2 as duas primeiras colunas de $\text{camera} R_{\text{tag}}$ e $t = \text{camera} t_{\text{tag}}$ [2]. Assim, a detecção dos cantos fornece correspondências suficientes para estimar H ; a partir disso, recupera-se uma estimativa inicial de pose, normalizando escala e impondo ortogonalidade em R .

Em nível prático, uma pose mais precisa é obtida ao refinar (R, t) por minimização do erro de reprojeção dos cantos observados. Denotando por \mathbf{u}_i os cantos medidos na imagem e por \mathbf{X}_i os cantos conhecidos no plano da tag, o refinamento pode ser formulado como

$$\min_{\text{camera} R_{\text{tag}} \in SO(3), \text{camera} t_{\text{tag}}} \sum_{i=1}^4 \left\| \mathbf{u}_i - \pi \left(K \left(\text{camera} R_{\text{tag}} \mathbf{X}_i + \text{camera} t_{\text{tag}} \right) \right) \right\|^2, \quad (2.12)$$

Figure 2.5: Geometria da retificação: cantos detectados \rightarrow estimação de homografia $H \rightarrow$ patch canônico para leitura do payload.

onde $\pi([x \ y \ z]^\top) = [x/z \ y/z]^\top$. Esta equação é a ponte entre percepção e controle: ao reduzir erro de reprojeção, reduz-se a incerteza na pose estimada, o que melhora diretamente a estabilidade do seguimento de trajetória em robótica móvel.

Para reforçar visualmente a ligação entre cantos, retificação e leitura de bits, recomenda-se inserir na Figura 2.5 um esquema com (a) a tag observada em perspectiva, (b) os quatro cantos detectados, e (c) o patch retificado (vista canônica) usado para amostragem do payload.

2.2.4 AprilTag 3 e confiabilidade sob generalização de layouts

O AprilTag 3 estende o sistema ao permitir layouts mais flexíveis, o que amplia possibilidades de desenho (incluindo regiões ignoradas e diferentes distribuições de bits), mas torna ainda mais importante controlar falsos positivos. O trabalho discute critérios de complexidade para garantir que os padrões gerados sejam improváveis em imagens naturais; entre eles, avalia-se a energia do modelo de Ising aplicada à imagem binária do layout [5]. Em termos discretos, uma escrita típica dessa métrica é

$$E = \sum_{i,j} \left(\mathbb{I}\{I(i,j) \neq I(i,j+1)\} + \mathbb{I}\{I(i,j) \neq I(i+1,j)\} \right), \quad (2.13)$$

que corresponde ao comprimento total de fronteiras preto-branco. A intuição é que padrões com estrutura local “muito regular” tendem a ocorrer com maior probabilidade em cenas artificiais/-naturais; logo, impor complexidade adequada reduz a chance de um padrão aleatório passar pelo processo de validação e ser aceito como tag [5].

3 Otmizacao e trajetoria

Apos um visualizaçao em que um camera se move em uma trajetoria que se paraece com incosaedro, durante a visualizaçao foi garantido que a todo momento a camera ve pelo menos duas tags;

A primeira etapa a realizar é fazer mapeamento do nosso ambiente. Sendo assim, definir a primeira tag referencia o calculadas tranformaçoes relativas para as tags subsequentes. dado que eu sei onde todas as tags estao localizada posição e orientação, eu posso localizar o meu robo no ambiente.

A problemantica é que o robo nao consegue ver todas as tags ao mesmo tempo e também, dado que o robo ve as tags em sequencia pequenos erros de detecçao da posicao do robo sao acumulados a medida em que o robo se move

Logo o mapeamento deve ser em caideia fazendo as transformada de referencial a partir de uma tag referencia, ou a tag mundo.

A camera é colocada como pose de referencia e recebe a pose relativa da tagA, primeira tag vista, ao mesmo tempo, é possivel ver a tagB, Assim atraves destas poses é possivel cacula a pose da tag B em relação a tagA:

$$T_{A \rightarrow B} = T_{Cam \rightarrow A}^{-1} * T_{Cam \rightarrow B} \quad (3.1)$$

- $T_{A \rightarrow B} = (R_A, t_A)$: Pose da Tag B em relação a A.
- $T_{Cam \rightarrow A} = (R_A, t_A)$: Pose da Tag A vista pela Câmera.
- $T_{Cam \rightarrow B} = (R_B, t_B)$: Pose da Tag B vista pela Câmera.
- R : Matriz de rotaçao
- t : Vetor de posição

A tranformação de inveresa é caculada da seguinte forma:

$$T^{-1} = (R^T, -R^T t) \quad (3.2)$$

dada que o as orientações das tags sao dadas em quartenion, é feita a conversao para matrizes de rotacao:

$$R = \begin{bmatrix} q_0^2 + q_1^2 - q_2^2 - q_3^2 & 2(q_1q_2 - q_0q_3) & 2(q_0q_2 + q_1q_3) \\ 2(q_1q_2 + q_0q_3) & q_0^2 - q_1^2 + q_2^2 - q_3^2 & 2(q_2q_3 - q_0q_1) \\ 2(q_1q_3 - q_0q_2) & 2(q_0q_1 + q_2q_3) & q_0^2 - q_1^2 - q_2^2 + q_3^2 \end{bmatrix} \quad (3.3)$$

Assim é possível calcular a rotação relativa (R_{rel}):

$$R_{rel} = R_A^T R_B \quad (3.4)$$

Dado que a detecção da tag também nos fornece a posição, assim também é possível calcular: a translação relativa (\mathbf{t}_{rel}):

$$\mathbf{t}_{rel} = R_A^T (\mathbf{t}_B - \mathbf{t}_A) \quad (3.5)$$

Assim

$$T_{A \rightarrow B} = (R_{rel}, \mathbf{t}_{rel}) \quad (3.6)$$

também é possível realizar o cálculo através de matrizes de transformação homogênea:

$$\begin{bmatrix} R_{rel} & \mathbf{t}_{rel} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} = \begin{bmatrix} R_A & \mathbf{t}_A \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix}^{-1} \begin{bmatrix} R_B & \mathbf{t}_B \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} \quad (3.7)$$

3.0.1 Subsubsection

4 Metodologia

5 Results

6 Discussion

7 Conclusion

References

- [1] P. Nourizadeh, M. Milford, and T. Fischer, *Teach and repeat navigation: A robust control approach*, 2024. arXiv: 2309.15405 [cs.R0]. [Online]. Available: <https://arxiv.org/abs/2309.15405>.
- [2] E. Olson, “AprilTag: A robust and flexible visual fiducial system,” in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, May 2011, pp. 3400–3407.
- [3] R. I. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, Second. Cambridge University Press, ISBN: 0521540518, 2004.
- [4] J. Wang and E. Olson, “AprilTag 2: Efficient and robust fiducial detection,” in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct. 2016.
- [5] M. Krogus, A. Haggemiller, and E. Olson, “Flexible layouts for fiducial tags,” in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct. 2019.