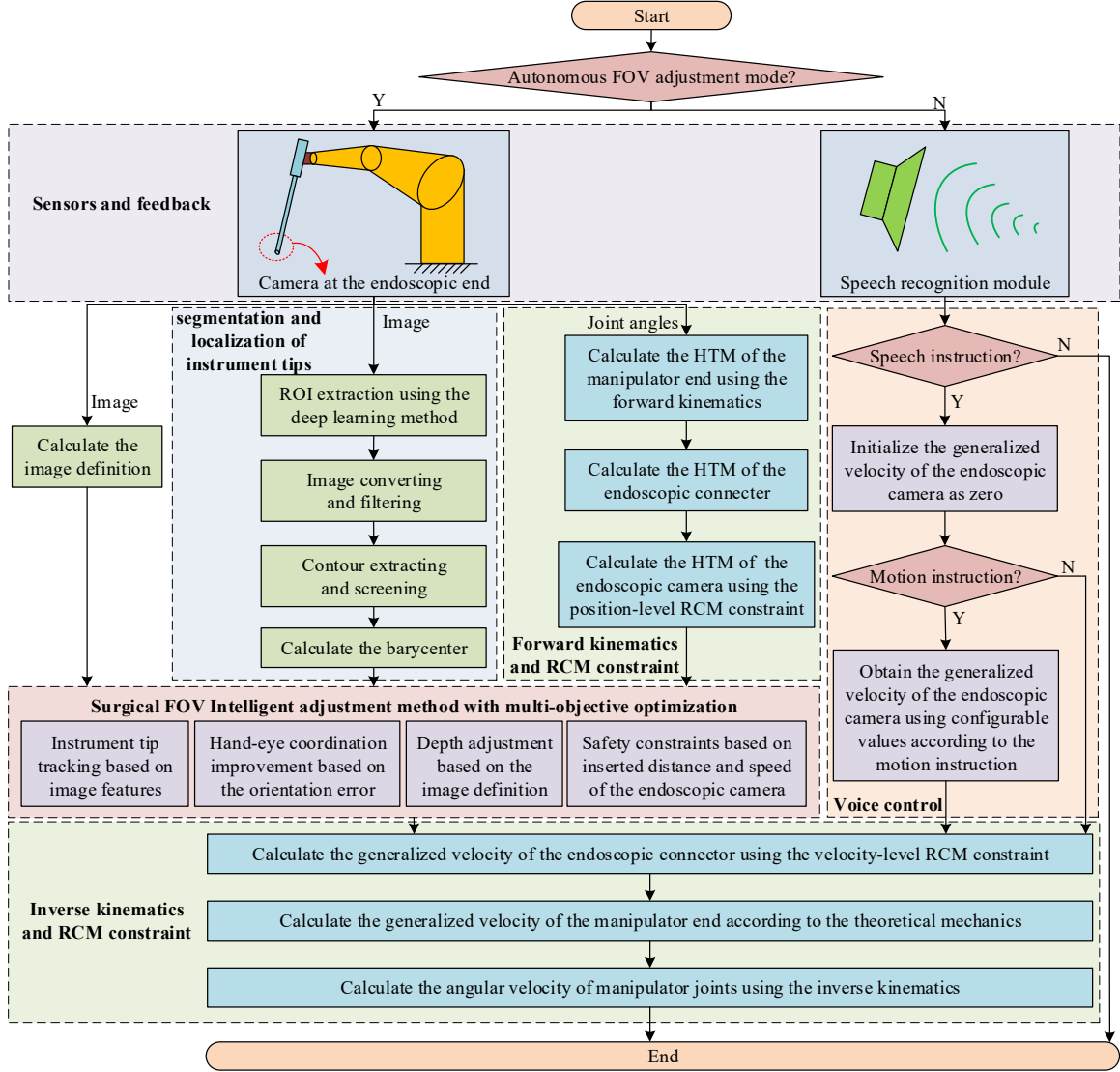## APPENDIX A

In order to better reflect the comprehensibility of the AEHR system and the control process, the calculation process of the entire algorithm is supplemented, as shown in **APPENDIX A.1**.



**APPENDIX A.1** The algorithm flowchart of the AEHR system
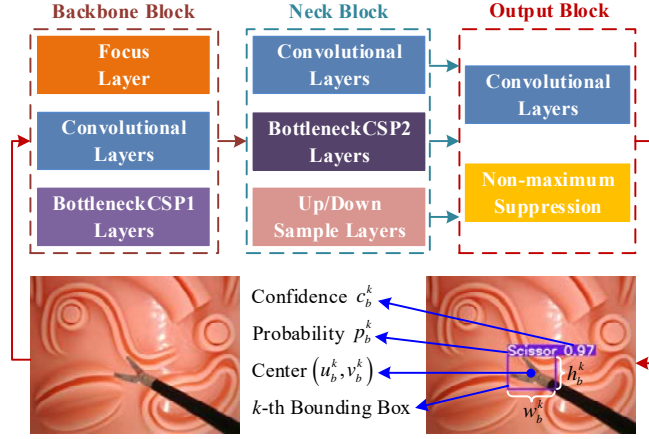
## APPENDIX B

The localization of surgical instrument tips can provide effective and important feedback information to intelligent FOV adjustment of AEHRs. In order to accurately extract the system tracking point, this paper adopts the DLM for segmentation and localization of instrument tips, which can reduce the localization error of instrument tips. Firstly, the region of interest (ROI) (i.e., the instrument tip region) is extracted in visual feedback using the YOLOv5 model. Then, the instrument tips are segmented, and the barycenter of instrument tips is extracted. Finally, the weighted average barycenter of all instrument tips is calculated as the system tracking point.

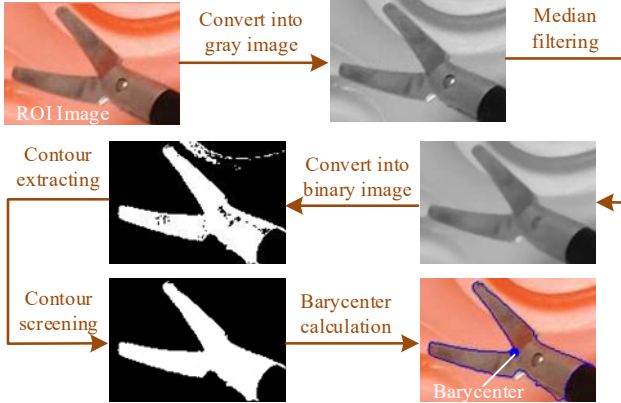**ROI extraction**: The YOLOv5 model has the ability of fast and accurate target region regression prediction, which can meet the real-time requirement of ROI extraction. It is assumed that $I_S$ is the $w_s \times h_s$ surgical image captured by the endoscopic camera, which contains the projection of $K$ instruments. The schematic diagram of ROI extraction is shown in **APPENDIX B.1**, and the prediction information (i.e., center, size, confidence, and class) of the $k$th instrument tip region is expressed as $\left(u_b^k, v_b^k, w_b^k, h_b^k, c_b^k, p_b^k\right)$.

**Segmentation and localization of instrument tips:** The segmentation and localization process of instrument tips is shown in **APPENDIX B.2**. It mainly includes the following steps: (1) crop out the $k$th ROI image $I_{TS}^k$ in $I_S$ according to $\left(u_b^k, v_b^k, w_b^k, h_b^k\right)$; (2) convert $I_{TS}^k$ into the gray image $I_{TG}^k$; (3) process $I_{TG}^k$ using median filter according to the filter size $k_s$

IEEE Transactions on Cybernetics

and obtain the filtered image $I_{\text{TGF}}^k$; (4) convert $I_{\text{TGF}}^k$ into the binary image $I_{\text{TB}}^k$ according to the threshold $[tr_{\min}, tr_{\max}]$; (5) extract $Q$ contours in $I_{\text{TB}}^k$; (6) calculate the area $a^q$ of $q$-th ($q = 1,2,...,Q$) contour $\boldsymbol{c}^q$ and the area ratio $\eta_{\text{area}} = a^q / \left( w_b^k \times h_b^k \right)$ of $\boldsymbol{c}^q$ in $I_{\text{TS}}^k$; (7) calculate the barycenter $\boldsymbol{s}_{\text{tip}}^k$ of $k$th instrument tip according to $\boldsymbol{c}^q$.



**APPENDIX B.1** The schematic diagram of ROI extraction



**APPENDIX B.2** The segmentation and localization process of instrument tips

The whole algorithm of the surgical instrument tips segmentation and localization is shown in **Algorithm B**.

---

**Algorithm** B Segmentation and localization of surgical instrument tips

---

**Input**: $I_S$ - source image captured by the endoscopic

**Output**: $\boldsymbol{s}_{\text{tips}}$ - system tracking point

$\left\{ \left( u_b^k, v_b^k, w_b^k, h_b^k, c_b^k, p_b^k \right) \mid k = 1,2,...,K \right\} \leftarrow$ input $I_S$ into the YOLOv5 model to predict ROIs

  **for** $k = 1$ to $K$ **do**

    $I_{\text{TS}}^k \leftarrow$ crop $I_S$ using $\left( u_b^k, v_b^k, w_b^k, h_b^k \right)$

    $I_{\text{TG}}^k \leftarrow$ convert $I_{\text{TS}}^k$ into a gray image

    $I_{\text{TGF}}^k \leftarrow I_{\text{TG}}^k$ using median filter, of which the size is $k_{\text{size}}$

    $I_{\text{TB}}^k \leftarrow$ convert $I_{\text{TGF}}^k$ into a binary image using $(tr_{\min}, tr_{\max})$

    $\left\{ \boldsymbol{c}^q \mid q = 1,2,...,Q \right\} \leftarrow$ extract contours of $I_{\text{TB}}^k$

    **for** $q = 1$ to $Q$ **do**

      $a^q \leftarrow$ calculate the area of $\boldsymbol{c}^q$

---

      calculate the area ratio of $\boldsymbol{c}^q$ in $I_{\text{TS}}^k$: $\eta_{\text{area}} = a^q / \left( w_b^k \times h_b^k \right)$

      **if** $\eta_{\text{area}} > ra_{\min}$ and $\eta_{\text{area}} < ra_{\max}$

        $I_{\text{TBS}}^k \leftarrow$ choose the region of $\boldsymbol{c}^q$ from $I_{\text{TB}}^k$

      **end if**

    **end for**

    $\boldsymbol{s}_{\text{tip}}^k \leftarrow$ (12) ~ (13) with $I_{\text{TBS}}^k$

  **end for**

  $\boldsymbol{s}_{\text{tips}} \leftarrow$ (14) with $\left\{ \boldsymbol{s}_{\text{tip}}^k \mid k = 1,2,...,K \right\}$

**return** $\boldsymbol{s}_{\text{tips}}$

---

## APPENDIX C

### Proof of Theorem

On the one hand, the proposed method decouples and controls the position and attitude of the endoscopic camera at the same time. On the other hand, it adjusts the endoscopic depth to optimize the image definition under the premise of meeting the position-level safety constraint.

**Theorem C1:**

According to the principle of pinhole imaging, it can be obtained:

$$\begin{bmatrix} u_{\text{tips}}^i \\ v_{\text{tips}}^i \\ 1 \end{bmatrix} = \frac{1}{{}^{\text{lap}}z_i} \begin{bmatrix} f_u & 0 & u_0 \\ 0 & f_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} {}^{\text{lap}}x_i \\ {}^{\text{lap}}y_i \\ {}^{\text{lap}}z_i \end{bmatrix} \quad (\text{C.1})$$

As shown in Fig.5, the direction vector of the endoscopic camera velocity is determined by:

$$\boldsymbol{v}_r^i = \frac{\boldsymbol{z}_{\text{lap}} \times \hat{\boldsymbol{p}}_{\text{tips}}^i \times \hat{\boldsymbol{p}}_{\text{tips}}^i}{\left\| \boldsymbol{z}_{\text{lap}} \times \hat{\boldsymbol{p}}_{\text{tips},i}^i \times \hat{\boldsymbol{p}}_{\text{tips}}^i \right\|} \quad (\text{C.2})$$

where, $\boldsymbol{z}_{\text{lap}} = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}^{\text{T}}$, $\hat{\boldsymbol{p}}_{\text{tips}}^i = \begin{bmatrix} x_{\text{tips}}^i / z_{\text{tips}}^i & y_{\text{tips},i}^i / z_{\text{tips}}^i & 1 \end{bmatrix}^{\text{T}}$, $\boldsymbol{p}_{\text{tips}}^i = [\, x_{\text{tips}}^i, y_{\text{tips}}^i, z_{\text{tips}}^i \,]^{\text{T}}$ is the 3D coordinate of the $i$th device tracking point.

Therefore, the linear velocity of the endoscopic camera can be expressed as:

$$\boldsymbol{v}_c^i = \begin{bmatrix} v_{cx}^i \\ v_{cy}^i \\ v_{cz}^i \end{bmatrix} = \upsilon_s^i \boldsymbol{v}_r^i = \upsilon_s^i \begin{bmatrix} x_{\text{tips}}^i z_{\text{tips}}^i / l \\ y_{\text{tips},i}^i z_{\text{tips}}^i / l \\ \left( x_{\text{tips}}^{i,2} + y_{\text{tips}}^{i,2} \right) / l \end{bmatrix} \quad (\text{C.3})$$

where, $l = \sqrt{\left( x_{\text{tips}}^{i\,2} + y_{\text{tips}}^{i\,2} \right)} \left\| \boldsymbol{p}_{\text{tips}}^i \right\|$, $\upsilon^i \geq 0$ is the magnitude of the camera velocity, $\boldsymbol{v}_r^i = \begin{bmatrix} v_{rx}^i, v_{ry}^i, v_{rz}^i \end{bmatrix}^{\text{T}}$.

When the endoscopic camera moves and the surgical instrument tip is stationary, the velocity of the surgical instrument tip relative to the endoscopic camera is:

$$\boldsymbol{v}_f^i = \begin{bmatrix} -v_{rx}^i - z_{\text{tips}}^i \omega_{ry}^i \\ -v_{ry}^i + z_{\text{tips}}^i \omega_{rx}^i \\ -v_{rz}^i - y_{\text{tips}}^i \omega_{rx}^i + x_{\text{tips}}^i \omega_{ry}^i \end{bmatrix} \quad (\text{C.4})$$

where, $\omega_{rx}^i = -v_{ry}^i / d_{in}$, $\omega_{ry}^i = v_{rx}^i / d_{in}$.

When both the endoscopic camera and the surgical instrument tips are moving, the linear velocity of the surgical

IEEE Transactions on Cybernetics

instrument tips relative to the endoscopic camera is:

$$\hat{\boldsymbol{v}}_f^i = \begin{bmatrix} \dot{x}_{tips}^i - v_{rx}^i - \dfrac{v_{rx}^i z_{tips}^i}{d_{in}} \\[2ex] \dot{y}_{tips}^i - v_{ry}^i - \dfrac{v_{ry}^i z_{tips}^i}{d_{in}} \\[2ex] \dot{z}_{tips}^i - v_{rz}^i + \dfrac{v_{ry}^i y_{tips}^i}{d_{in}} + \dfrac{v_{rx}^i x_{tips}^i}{d_{in}} \end{bmatrix}$$

$$= \begin{bmatrix} \dot{x}_{tips}^i - \dfrac{\upsilon^i}{l} x_{tips,i}^i z_{tips,i} - \dfrac{\upsilon^i}{d_{in}l} x_{tips,i}^i z_{tips}^{i2} \\[2ex] \dot{y}_{tips}^i - \dfrac{\upsilon^i}{l} y_{tips}^i z_{tips}^i - \dfrac{\upsilon^i}{d_{in}l} y_{tips}^i z_{tips}^{i,2} \\[2ex] \dot{z}_{tips}^i + \dfrac{\upsilon^i}{l}\left(x_{tips}^{i,2} + y_{tips}^{i,2}\right) + \dfrac{\upsilon^i}{d_{in}l}\left(x_{tips}^{i,2} + y_{tips}^{i,2}\right) z_{tips}^i \end{bmatrix}$$

$$(C.5)$$

Differentiate (C.1) and substitute (C.5) into the result of the derivation to get:

$$\dot{u}_{tips}^i = \frac{f_u}{z_{tips}^i} v_{fx}^i - \frac{f_u x_{tips}^i}{z_{tips}^{i2}} v_{fz}^i$$

$$= \frac{f_u}{z_i^{i2}}\Big[ \dot{x}_{tips}^i z_{tips}^i - x_{tips}^i \dot{z}_{tips}^i$$

$$-\frac{\upsilon^i}{l} x_{tips}^i \left( x_{tips}^{i2} + y_{tips}^{i2} + z_{tips}^{i2} \right)$$

$$-\frac{\upsilon^i}{d_{in}l} x_{tips}^i z_{tips}^i \left( x_{tips}^{i2} + y_{tips}^{i2} + z_{tips}^{i2} \right) \Big] \quad (C.6)$$

$$\dot{v}_{tips}^i = \frac{f_v}{z_{tips}^i} v_{fy}^i - \frac{f_v y_{tips}^i}{z_i^2} v_{fz}^i$$

$$= \frac{f_v}{z_{tips}^{i2}}\Big[ \dot{y}_{tips}^i z_{tips}^i - y_{tips}^i \dot{z}_{tips}^i$$

$$-\frac{\upsilon^i}{l} y_{tips}^i \left( x_{tips}^{i2} + y_{tips}^{i2} + z_{tips}^{i2} \right)$$

$$-\frac{\upsilon^i}{d_{in}l} y_{tips}^i z_{tips}^i \left( x_{tips}^{i2} + y_{tips}^{i2} + z_{tips}^{i2} \right) \Big] \quad (C.7)$$

The established Lyapunov function is the distance between the pixel coordinates of the actual instrument tips and the center of the camera imaging image, that is,

$$L = \left( u_{tips}^i - u_0 \right)^2 + \left( v_{tips}^i - v_0 \right)^2 \quad (C.8)$$

Taking the derivative of (C.8), we can get:

$$\dot{L} = 2\left( u_{tips}^i - u_0 \right)\dot{u}_{tips}^i + 2\left( u_{tips}^i - v_0 \right)\dot{v}_{tips}^i \quad (C.9)$$

Substituting (C.1), (C.6), and (C.7) into (C.9), yield:

$$\dot{L} = \frac{f_u^2 x_{tips}^{i2}}{z_{tips}^{i3}}\left[ \frac{z_{tips}^i}{x_{tips}^i}\dot{x}_{tips}^i - \dot{z}_{tips}^i - \left(1 + \frac{z_{tips}^i}{d_{in}}\right)\upsilon^i \frac{\|\boldsymbol{p}_{tips}^i\|}{\sqrt{x_{tips}^{i2} + y_{tips}^{i2}}} \right]$$

$$+ \frac{f_v^2 y_{tips}^{i2}}{z_{tips}^{i3}}\left[ \frac{z_{tips}^i}{y_{tips}^i}\dot{y}_{tips}^i - \dot{z}_{tips}^i - \left(1 + \frac{z_{tips}^i}{d_{in}}\right)\upsilon^i \frac{\|\boldsymbol{p}_{tips}^i\|}{\sqrt{x_{tips}^{i2} + y_{tips}^{i2}}} \right]$$

$$(C.10)$$

When $d_{in} \to +\infty$, (10) can be further simplified as:

$$\dot{L} = \frac{f_u^2 x_{tips}^{i2}}{z_{tips}^{i3}}\left[ \frac{z_{tips}^i}{x_{tips}^i}\dot{x}_{tips}^i - \dot{z}_{tips}^i - \frac{\|\boldsymbol{p}_{tips}^i\|}{\sqrt{x_{tips}^{i2} + y_{tips}^{i2}}}\upsilon^i \right] +$$

$$\frac{f_v^2 y_{tips}^{i2}}{z_{tips}^{i3}}\left[ \frac{z_{tips}^i}{y_{tips}^i}\dot{y}_{tips}^i - \dot{z}_{tips}^i - \frac{\|\boldsymbol{p}_{tips}^i\|}{\sqrt{x_{tips}^{i2} + y_{tips}^{i2}}}\upsilon^i \right] \quad (C.11)$$

$$\leq \frac{f_u^2 x_{tips}^{i2}}{z_{tips}^{i3}} A + \frac{f_v^2 y_{tips}^{i2}}{z_{tips}^{i3}} B$$

with

$$A = \frac{z_{tips}^i}{x_{tips}^i}\dot{x}_{tips}^i - \dot{z}_{tips}^i - \frac{\|\boldsymbol{p}_{tips}^i\|}{\sqrt{x_{tips}^{i2} + y_{tips}^{i2}}}\upsilon^i$$

$$B = \frac{z_{tips}^i}{y_{tips}^i}\dot{y}_{tips}^i - \dot{z}_{tips}^i - \frac{\|\boldsymbol{p}_{tips}^i\|}{\sqrt{x_{tips}^{i2} + y_{tips}^{i2}}}\upsilon^i$$

$$(C.12)$$

Since $z_{tips}^i > 0$, (12) has the following characteristics:

$$A \leq \max\left\{ \frac{z_{tips}^i}{\|x_{tips}^i\|}\upsilon^i{}_s, \upsilon^i \right\} - \frac{\|\boldsymbol{p}_{tips}^i\|}{\sqrt{x_{tips}^{i2} + y_{tips}^{i2}}}\upsilon^i$$

$$B \leq \max\left\{ \frac{z_{tips}^i}{\|y_{tips}^i\|}\upsilon^i, \upsilon^i \right\} - \frac{\|\boldsymbol{p}_{tips}^i\|}{\sqrt{x_{tips}^{i2} + y_{tips}^{i2}}}\upsilon^i$$

$$(C.13)$$

Further,

$$A: \begin{cases} A \leq \upsilon^i - \dfrac{\|\boldsymbol{p}_{tips}^i\|}{\sqrt{x_{tips}^{i2} + y_{tips}^{i2}}}\upsilon^i \leq 0, \text{ if } \dfrac{z_{tips}^i}{\|x_{tips}^i\|} \leq 1 \\[3ex] A \leq \dfrac{z_{tips}^i}{\|x_{tips}^i\|}\upsilon^i - \sqrt{\dfrac{x_{tips}^{i2} + z_{tips}^{i2}}{x_{tips}^{i2}}}\upsilon^i \leq 0, \text{ else} \end{cases}$$

$$B: \begin{cases} B \leq \upsilon^i - \dfrac{\|\boldsymbol{p}_{tips}^i\|}{\sqrt{x_{tips}^{i2} + y_{tips}^{i2}}}\upsilon^i \leq 0, \text{ if } \dfrac{z_{tips}^i}{\|y_{tips}^i\|} \leq 1 \\[3ex] B \leq \dfrac{z_{tips}^i}{\|y_{tips}^i\|}\upsilon^i - \sqrt{\dfrac{y_{tips}^{i2} + z_{tips}^{i2}}{y_{tips}^{i2}}}\upsilon^i \leq 0, \text{ else} \end{cases}$$

$$(C.14)$$

Based on (C.14) and (C.11), $\dot{L} \leq 0$ can be obtained, so the distance error between the pixel coordinates of the actual instrument end and the center of the camera image is always convergent.

**Theorem C2:**

According to the hand-eye coordination model, $\exists \beta_1 \geq 0$, when $\forall \beta \geq \beta_1$, $0 \leq {}^{con}\omega_{coor} \leq \omega_{max}$.

$$\dot{\beta} = -\eta_{coor}{}^{con}\omega_{coor} \leq 0 \quad (C.15)$$

Empathy, $\exists \beta_2 \leq 0$, when $\forall \beta \leq \beta_2$, $-\omega_{max} \leq {}^{con}\omega_{coor}^i \leq 0$.

$$\dot{\beta} = -\eta_{coor}{}^{con}\omega_{coor} \geq 0 \quad (C.16)$$

**Theorem C3:**

IEEE Transactions on Cybernetics

Let $\eta_{\text{def}} \geq \eta_{\text{track}}$. According to the optimization model considering image definition, $\exists |\eta_q| \geq \eta_0$, when $\forall d_{obj} < d_{foc}$, $\lambda_{\text{def}} = v_{\max}$.

$$\dot{d}_{\text{obj}} = \eta_{\text{track}} {}^{lap}\boldsymbol{v}_{\text{track}} \cdot \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}^{\text{T}} + \eta_{\text{def}} \lambda_{\text{def}}$$
$$\geq \left( -\eta_{\text{track}} + \eta_{\text{def}} \right) v_{\max} \qquad \text{(C.17)}$$
$$\geq 0$$

Empathy, $\exists |\eta_q| \geq \eta_0$, when $\forall d_{obj} > d_{foc}$, $\lambda_{\text{def}} = -v_{\max}$.

$$\dot{d}_{\text{obj}} = \eta_{\text{track}} {}^{lap}\boldsymbol{v}_{\text{track}} \cdot \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}^{\text{T}} + \eta_{\text{def}} \lambda_{\text{def}}$$
$$\leq \left( \eta_{\text{track}} - \eta_{\text{def}} \right) v_{\max} \qquad \text{(C.18)}$$
$$\leq 0$$

**Theorem C4:**

Let $\eta_{\text{safe}} \geq \eta_{\text{track}}$. According to the position-level safety constraint, $\exists d_{in1} \in \left[ d_{\text{alart,min}}, d_{\text{safe,min}} \right]$, when $\forall d_{in} \in \left[ d_{\text{alart,min}}, d_{in1} \right]$, $\lambda_{safe} = v_{\max}$.

$$\dot{d}_{in} = \eta_{\text{track}} {}^{con}\boldsymbol{v}_{\text{track}} \cdot \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}^{\text{T}} + \eta_{\text{safe}} \lambda_{safe}$$
$$\geq \left( -\eta_{\text{track}} + \eta_{\text{safe}} \right) v_{\max} \qquad \text{(C.19)}$$
$$\geq 0$$

Empathy, $\exists d_{in2} \in \left[ d_{\text{safe,max}}, d_{\text{alart,max}} \right]$, when $\forall d_{in} \in \left[ d_{in2}, d_{\text{alart,max}} \right]$, $\lambda_{safe} = -v_{\max}$.

$$\dot{d}_{in} = \eta_{\text{track}} {}^{con}\boldsymbol{v}_{\text{track}} \cdot \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}^{\text{T}} + \eta_{\text{safe}} \lambda_{safe}$$
$$\leq \left( \eta_{\text{track}} - \eta_{\text{safe}} \right) v_{\max} \qquad \text{(C.20)}$$
$$\leq 0$$

## APPENDIX D

In the experimental system, the D-H parameters of the 6-DOF robot are shown in **APPENDIX D.1**. During instrument tips identification, the architecture of surgical instrument detection with YOLOv5 model is shown in **APPENDIX D.2.**

**APPENDIX D.1** The D-H parameter of the robot

| Link $i$ | $\alpha_i$ (rad) | $a_i$ (mm) | $d_i$ (mm) | offset(rad) | $\theta_i$ (rad) |
|---|---|---|---|---|---|
| 1 | $\pi/2$ | 0 | 199.34 | $-\pi$ | $\theta_1$ |
| 2 | 0 | -250.00 | 0 | 0 | $\theta_2$ |
| 3 | 0 | -250.00 | 0 | 0 | $\theta_3$ |
| 4 | $\pi/2$ | 0 | 109.10 | 0 | $\theta_4$ |
| 5 | $-\pi/2$ | 0 | 108.00 | 0 | $\theta_5$ |
| 6 | 0 | 0 | 75.86 | 0 | $\theta_6$ |

**APPENDIX D.2** Architecture of surgical instrument detection with YOLOv5 model

| Structure name | Filter size | Input shape | Output shape |
|---|---|---|---|
| Adaptive resize | - | 640×480×3 | 608×608×3 |
| Focus | 32 | 608×608×3 | 304×304×32 |
| BatchNormalization+ Convolution | 64 | 304×304×32 | 152×152×64 |
| BottleneckCSP | 64 | 152×152×64 | 152×152×64 |
| BatchNormalization+ Convolution | 128 | 152×152×64 | 76×76×128 |
| BottleneckCSP×3 | 128 | 76×76×128 | 76×76×128 |
| BatchNormalization+ Convolution | 256 | 76×76×128 | 38×38×256 |
| BottleneckCSP×3 | 256 | 38×38×256 | 38×38×256 |
| BatchNormalization+ Convolution | 512 | 38×38×256 | 19×19×512 |
| Neck | - | 76×76×128<br>38×38×256<br>19×19×512 | 76×76×255<br>38×38×255<br>19×19×255 |