



Código	FPI-009
Objeto	Guía de elaboración de Informe final de proyecto
Usuario	Director de proyecto de investigación
Autor	Secretaría de Ciencia y Tecnología de la UNLaM
Versión	5
Vigencia	03/9/2019

**Departamento:
Ingeniería e Investigaciones Tecnológicas**

**Programa de acreditación:
PROINCE**

Programa de Investigación¹:

**Código del Proyecto:
C218**

**Título del proyecto:
Obtención de data set a partir de imágenes capturadas del estacionamiento de UNLaM para su
utilización en sistema experto de reconocimiento de imágenes**

**Director:
Ing. Fernando I. Szklanny**

**Codirector:
Lic. Carlos Eduardo Maidana**

**Integrantes:
Ing. Edgardo Gho
Ing. Martín Ferreyra Birón
Ing. Carlos Alberto Rodríguez**

**Resolución Rectoral de acreditación:
N°354/19**

**Fecha de inicio:
01/01/2019**

**Fecha de finalización:
30/04/2021**

¹ Los Programas de Investigación de la UNLaM están acreditados con resolución rectoral, según lo indica la Resolución HCS N° 014/15 sobre **Lineamientos generales para el establecimiento, desarrollo y gestión de Programas de Investigación a desarrollarse en la Universidad Nacional de La Matanza**. Consultar en el departamento académico correspondiente la inscripción del proyecto en un Programa acreditado.



Código	FPI-009
Objeto	Guía de elaboración de Informe final de proyecto
Usuario	Director de proyecto de investigación
Autor	Secretaría de Ciencia y Tecnología de la UNLaM
Versión	5
Vigencia	03/9/2019

A. Desarrollo del proyecto (adjuntar el protocolo)

A.1. Tal como se hizo referencia en el informe de avance el presente proyecto tuvo durante su primer año de desarrollo una serie de inconvenientes, que obligaron a los investigadores a replantear varias veces el objetivo y el enfoque del mismo.

El proyecto se inició original y teóricamente en el mes de febrero de 2019. Al comenzar a trabajar más profundamente en el tema, surgió la necesidad de algún equipamiento adicional al que posee el grupo de investigación, así como la mejora en las prestaciones de las computadoras disponibles.

Una vez adquirido el equipamiento adicional se requirió confirmación de la autorización para instalar la cámara de video en la playa de estacionamiento.

Debido a diversos inconvenientes de tipo edilicio en la zona en la que se debía instalar la mencionada cámara, se buscó, como alternativa, replantear el proyecto, cuya base era la recolección de imágenes de un conjunto conocido de patentes vehiculares para la creación de una base de datos.

Se planteó como alternativa temporaria la adquisición de datos en la vía pública por medio de una cámara tipo "Dashcam". Esta solución requiere la toma de varias horas de imágenes en video y la posterior extracción manual de cada una de las patentes, para luego seleccionar las que fuesen lo suficientemente legibles como para ser incorporadas al acervo.

A comienzos del año 2020 se continuó trabajando con la presente investigación, pero todavía con las restricciones internas respecto a la captura de imágenes dentro de la universidad ya mencionadas. Estas restricciones, que se debieron a razones ajenas a este grupo de investigación, de alguna forma se encontraban subsanadas y, solamente, restaba esperar a instalar las cámaras de video en el estacionamiento de la universidad y a tener un flujo importante de autos en el estacionamiento, pero, por esos meses el Aislamiento Social Preventivo y Obligatorio comenzó a regir, la universidad cerró las puertas y el flujo de automóviles para capturar imágenes de automóviles y patentes fue nulo en la universidad, y en los primeros meses de la pandemia, también en gran parte del AMBA.

Esta situación, junto a las obligaciones adicionales sumadas producto del aislamiento perjudicó sustancialmente el progreso de esta investigación, sin embargo, a pesar de toda esta adversidad se analizó cuál sería la mejor manera de continuar con la misma. Nuevamente se debió cambiar el rumbo y adaptarse a las circunstancias, para subsanar ello, se comenzó a confeccionar un programa para segmentar de forma manual las patentes, para así comenzar a armar el acervo de datos. De esta forma los integrantes del grupo instalaron cámaras en los parabrisas de sus automóviles y se comenzó a grabar videos. Por lo tanto la forma en la que se continuó con la presente investigación fue en un principio tomar el programa en cuestión y mejorarlo para comenzar a crear el acervo. Se debe tener en cuenta que al usar las imágenes del estacionamiento de la universidad, se hubiera podido corroborar la chapa patente contrastando la base de datos de la universidad y el código del TAG instalado en cada vehículo, información que se posee, ya que el grupo fue el encargado del desarrollo del sistema de estacionamiento de la universidad. Al filmar autos que no existen en esa base de datos se deben corroborar las imágenes manualmente.



Código	FPI-009
Objeto	Guía de elaboración de Informe final de proyecto
Usuario	Director de proyecto de investigación
Autor	Secretaría de Ciencia y Tecnología de la UNLaM
Versión	5
Vigencia	03/9/2019

Una de las primeras hipótesis planteadas para el correcto reconocimiento de patentes es que, lo primero que se debe buscar en la imagen, es la posición del automóvil en la imagen, o dicho de otra manera, lo primero que se debe segmentar de la imagen es el automóvil. Esto debe ser así ya que una patente, conceptualmente, es solo un cartel, y realizar la búsqueda de este cartel en una imagen no es eficiente e incluso puede ser contraproducente debido a los falsos positivos que pueden crearse, esto aún es peor cuando las imágenes son realmente muy grandes. Por lo tanto el programa anteriormente mencionado, se adaptó no solo para segmentar patentes y caracteres sino también para crear acervos de las tres clases anteriormente mencionadas. Toda la información segmentada que conforma este acervo es guardado en una carpeta acompañado de un archivo XML que "etiqueta" cada una de las respectivas imágenes dándole así un contexto a los datos.

Es importante destacar que los distintos puntos que se buscaban alcanzar, como la creación de un acervo de datos y la detección de patentes en automóviles no se encuentra muy lejana de llevarse a cabo por completo: El acervo de datos puede crecer rápidamente debido a las herramientas desarrolladas, y si fue posible detectar automóviles con la red neuronal YOLO (You Only Look Once), también será posible entrenar a la red neuronal para detectar patentes en los autos, tanto patentes vigentes desde 1994 a 2016 en la República Argentina como también las que respetan el formato MERCOSUR vigente desde abril de 2016. Probablemente esto implique un desafío de implementación, pero el objetivo probablemente pueda ser alcanzado. Otro punto a investigar es el comportamiento de redes neuronales similares para realizar la tarea de detección de objetos en tiempo real y poder compararla con YOLO, y por último resta investigar si la utilización de la red YOLO es útil en la detección de caracteres en tiempo real, o existe alguna alternativa para un problema que en teoría ya no es un desafío para el campo de la inteligencia artificial debido a las diversas redes neuronales que a lo largo de los años se desarrollaron para cumplir con este objetivo.



Código	FPI-009
Objeto	Guía de elaboración de Informe final de proyecto
Usuario	Director de proyecto de investigación
Autor	Secretaría de Ciencia y Tecnología de la UNLaM
Versión	5
Vigencia	03/9/2019

B. Principales resultados de la investigación

B.1. Publicaciones en revistas (informar cada producción por separado)

Artículo 1:	
Autores:	Martín Ferreyra Birón, Fernando I. Szklanny, Alberto Miguens
Título del artículo:	Desarrollo de un sistema de asistencia visual para un manipulador automático programable
N° de fascículo:	1
N° de Volumen:	5
Revista:	<i>ReDDI</i>
Año:	2020
Institución editora de la revista	<i>UNLaM</i>
País de procedencia de institución editora	<i>Argentina</i>
Arbitraje	NO
ISSN:	2525-1333-Redd
URL de descarga del artículo	https://reddi.unlam.edu.ar/index.php/ReDDi/article/view/120/222
N° DOI	

B.2. Libros

Libro 1	
Autores	
Título del Libro	
Año	
Editorial	
Lugar de impresión	
Arbitraje	Elija un elemento.
ISBN:	
URL de descarga del libro	
N° DOI	



Código	FPI-009
Objeto	Guía de elaboración de Informe final de proyecto
Usuario	Director de proyecto de investigación
Autor	Secretaría de Ciencia y Tecnología de la UNLAM
Versión	5
Vigencia	03/9/2019

B.3. Capítulos de libros

Autores	
Título del Capítulo	
Título del Libro	
Año	
Editores del libro/Compiladores	
Lugar de impresión	
Arbitraje	Elija un elemento.
ISBN:	
URL de descarga del capítulo	
N° DOI	

B.4. Trabajos presentados a congresos y/o seminarios

Autores	<i>Ing. Martín Ferreyra Biron, Ing. Alberto Miguens, Ing. Fernando Szklanny</i>
Título	<i>Desarrollo de un sistema de asistencia visual...</i>
Año	<i>2019</i>
Evento	<i>CONAIISI 2019</i>
Lugar de realización	<i>UNLAM</i>
Fecha de presentación de la ponencia	<i>15/09/2019</i>
Entidad que organiza	<i>UNLAM</i>
URL de descarga del trabajo (especificar solo si es la descarga del trabajo; formatos pdf, e-pub, etc.)	<i>https://conaiisi2019.unlam.edu.ar/</i>

B.5. Otras publicaciones

Autores	
Año	
Título	
Medio de Publicación	



Código	FPI-009
Objeto	Guía de elaboración de Informe final de proyecto
Usuario	Director de proyecto de investigación
Autor	Secretaría de Ciencia y Tecnología de la UNLaM
Versión	5
Vigencia	03/9/2019

C. Otros resultados. Indicar aquellos resultados pasibles de ser protegidos a través de instrumentos de propiedad intelectual, como patentes, derechos de autor, de rechos de obtentor, etc. y desarrollos que no pueden ser protegidos por instrumentos de propiedad intelectual, como las tecnologías organizacionales y otros. Complete un cuadro por cada uno de estos dos tipos de productos.

C.1. Títulos de propiedad intelectual. Indicar: Tipo (marcas, patentes, modelos y diseños, la transferencia tecnológica) de desarrollo o producto, Titular, Fecha de solicitud, Fecha de otorgamiento

Tipo	Titular	Fecha de Solicitud	Fecha de Emisión

C.2. Otros desarrollos no pasibles de ser protegidos por títulos de propiedad intelectual. Indicar: Producto y Descripción.

Producto	Descripción

D. Formación de recursos humanos. Trabajos finales de graduación, tesis de grado y posgrado. Completar un cuadro por cada uno de los trabajos generados en el marco del proyecto.

D.1. Tesis de grado

Director (apellido y nombre)	Autor (apellido y nombre)	Institución	Calificación	Fecha /En curso	Título de la tesis

D.2 Trabajo Final de Especialización

Director (apellido y nombre)	Autor (apellido y nombre)	Institución	Calificación	Fecha /En curso	Título del Trabajo Final

D.2. Tesis de posgrado: Maestría

Director (apellido y nombre)	Tesista (apellido y nombre)	Institución	Calificación	Fecha /En curso	Título de la tesis



Código	FPI-009
Objeto	Guía de elaboración de Informe final de proyecto
Usuario	Director de proyecto de investigación
Autor	Secretaría de Ciencia y Tecnología de la UNLaM
Versión	5
Vigencia	03/9/2019

D.3. Tesis de posgrado: Doctorado

Director (apellido y nombre)	Tesista (apellido y nombre)	Institución	Calificación	Fecha /En curso	Título de la tesis

D.4. Trabajos de Posdoctorado

Director (apellido y nombre)	Posdoctorando (apellido y nombre)	Institución	Calificación	Fecha /En curso	Título del trabajo	Publicación

E. Otros recursos humanos en formación: estudiantes/ investigadores (grado/posgrado/ posdoctorado)

Apellido y nombre del Recurso Humano	Tipo	Institución	Período (desde/hasta)	Actividad asignada ²

F. Vinculación³: Indicar conformación de redes, intercambio científico, etc. con otros grupos de investigación; con el ámbito productivo o con entidades públicas. Desarrolle en no más de dos (2) páginas.

G. Otra información. Incluir toda otra información que se considere pertinente.

² Descripción de la/s actividad/es a cargo (máximo 30 palabras)

³Entendemos por acciones de “vinculación” aquellas que tienen por objetivo dar respuesta a problemas, generando la creación de productos o servicios innovadores y confeccionados “a medida” de sus contrapartes.



Código	FPI-009
Objeto	Guía de elaboración de Informe final de proyecto
Usuario	Director de proyecto de investigación
Autor	Secretaría de Ciencia y Tecnología de la UNLaM
Versión	5
Vigencia	03/9/2019

H. Cuerpo de anexos:

- Anexo I: Descripción detallada del proyecto.
- Anexo II: Copia de cada uno de los trabajos mencionados en los puntos B, C y D, y certificaciones cuando corresponda.⁴
- Anexo III:
 - FPI-013: Evaluación de alumnos integrantes. (si corresponde)
 - FPI-014: Comprobante de liquidación y rendición de viáticos. (si corresponde)
 - FPI-015: Rendición de gastos del proyecto de investigación acompañado de las hojas foliadas con los comprobantes de gastos.
 - FPI-035: Formulario de reasignación de fondos en Presupuesto.
- Anexo IV: Alta patrimonial de los bienes adquiridos con presupuesto del proyecto (FPI 017)
- Nota justificando baja de integrantes del equipo de investigación.

Lic. Carlos Eduardo Maidana
Codirector

Lugar y fecha: San Justo 30 de abril de 2021

⁴En caso de libros, podrá presentarse una fotocopia de la primera hoja significativa o su equivalente y el índice.

ANEXO I

Descripción detallada del proyecto

A comienzos del año 2020 se continuó trabajando con la presente investigación, pero todavía con las restricciones internas respecto a la captura de imágenes dentro de la universidad ya mencionadas en el informe anterior. Estas restricciones, que se debieron a razones ajenas a este grupo de investigación, de alguna forma se encontraban subsanadas y solamente, restaba esperar a instalar las cámaras de video en el estacionamiento de la universidad y poder tener un flujo importante de autos en el estacionamiento, pero, por esos meses el Aislamiento Social Preventivo y Obligatorio comenzó a regir, la universidad cerró las puertas y el flujo de automóviles para capturar imágenes de automóviles y patentes fue prácticamente nulo en la universidad y en los primeros meses de la pandemia en gran parte del AMBA. Esta situación junto a las obligaciones adicionales sumadas producto del aislamiento perjudicó sustancialmente el progreso de esta investigación, sin embargo, a pesar de toda esta situación se analizó cuál sería la mejor manera de continuar con la misma. A fines del año 2019, ya con la negativa de no poder obtener imágenes de automóviles en la universidad, se comenzó a confeccionar un programa para segmentar de forma manual las patentes, para así comenzar a armar el acervo de datos, como también se comenzó a grabar videos con una cámara instalada en los vehículos de los integrantes. Por lo tanto la forma en la que se continuó con la presente investigación fue en un principio tomar el programa en cuestión y mejorarlo para comenzar a crear el acervo.

Una de las primeras hipótesis planteadas para el correcto reconocimiento de patentes es que, lo primero que se debe buscar en la imagen, es la posición del automóvil en la imagen, o dicho de otra manera, lo primero que se debe segmentar de la imagen es el automóvil. Esto debe ser así ya que una patente, conceptualmente, es solo un cartel, y realizar la búsqueda de este cartel en una imagen no es eficiente e incluso puede ser contraproducente debido a los falsos positivos que pueden crearse, esto aún es peor cuando las imágenes son realmente muy grandes. Por lo tanto el programa anteriormente mencionado, se adaptó no solo para segmentar patentes y caracteres sino también para crear acervos de las tres clases anteriormente mencionadas. Toda la información segmentada que conforma este acervo es guardado en una carpeta acompañado de un archivo XML que "etiqueta" cada una de las respectivas imágenes dándole así un contexto a los datos.

A continuación se puede observar la pantalla principal del programa confeccionado el cual permite elegir el video principal del que se obtendrán las imágenes para así segmentar un auto o una patente (Fig. 1) para luego continuar con la siguiente pantalla que permite seleccionar el tipo de dato que fue segmentado para finalmente poder guardarlo como se muestra en la Fig. 2.

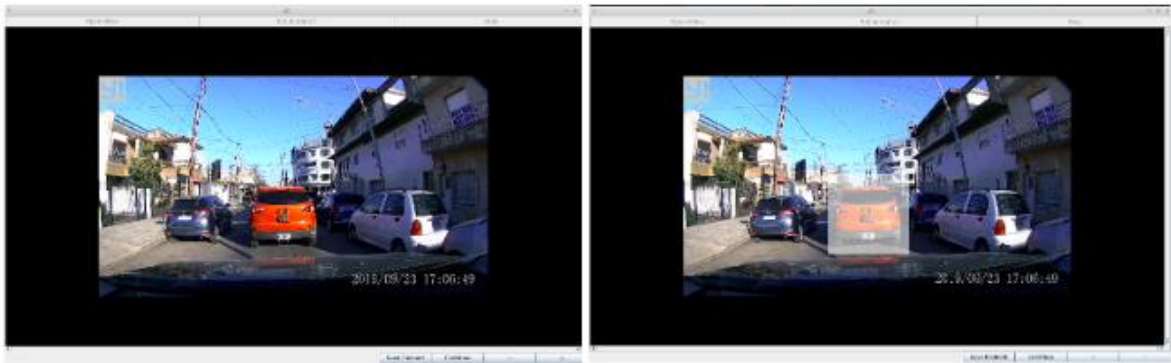
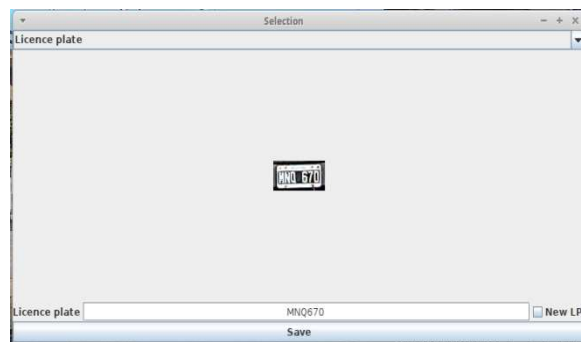


Fig. 1: Pantalla principal del programa desarrollado para segmentar autos, patentes y caracteres de patentes



A



B



C

Fig 2. A. Pantalla diseñada para anexar una imagen de un automóvil al acervo. B. Pantalla diseñada para anexar una patente al acervo. C. Pantalla para anexar los caracteres segmentados de las patentes.

Ya teniendo en nuestro poder este programa, nos enfocamos en parte, a segmentar manualmente autos, patentes y caracteres, empezando a cumplir con nuestro primer objetivo en la investigación que es

confeccionar un acervo de datos de imágenes de patentes. En este caso el acervo aumentó en clases debido a que segmentamos automóviles, patentes y caracteres de las mismas patentes. Habiéndose comenzado a lograr este objetivo también fue imperioso lograr visualizar de manera acertada este acervo, para así poder manipularlo, curar las imágenes que se obtuvieron si fuera necesario y poder administrarlo de alguna forma. Para esto se confeccionó otro programa el cual toma el archivo XML y muestra cada uno de los elementos que posee el acervo, con la posibilidad de modificar la imagen, o cualquier dato asociado al archivo XML. A continuación se puede observar la Fig 3. la cual muestra como se observan las tres clases o categorías en el programa

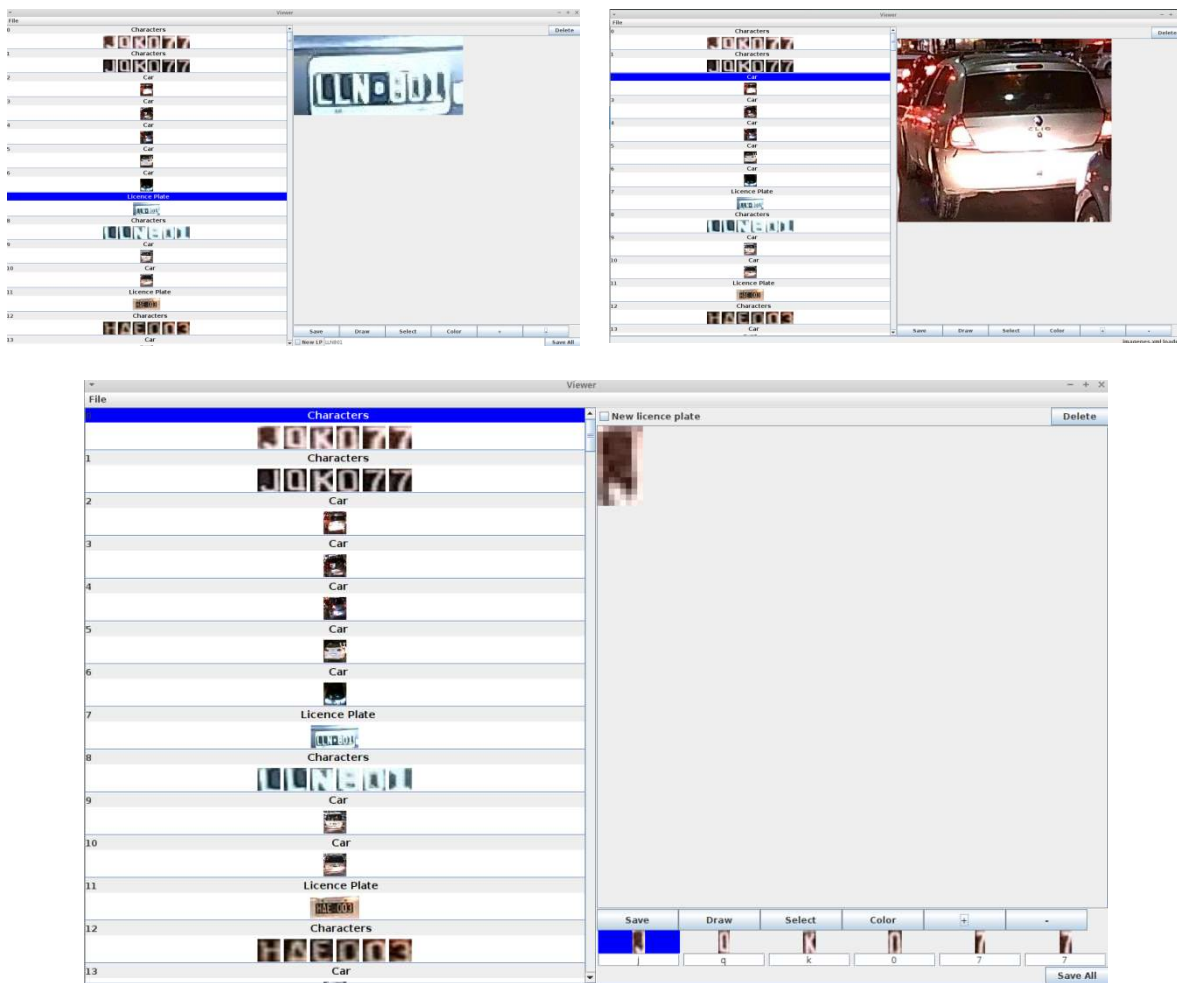


Fig. 3: Vista de la selección de un ejemplo de cada una de las clases segmentadas

Finalmente el programa permite guardar las imágenes de cada una de las clases segmentadas, (autos, caracteres, patentes o un carácter específico) y ejecutar por cada uno de los elementos guardados un comando particular, si es que se lo desea (por ejemplo si se desea redimensionar las imágenes se puede eje-

cutar por cada una de las imágenes el programa convert). La Fig. 4 muestra la ventana que permite esta acción.

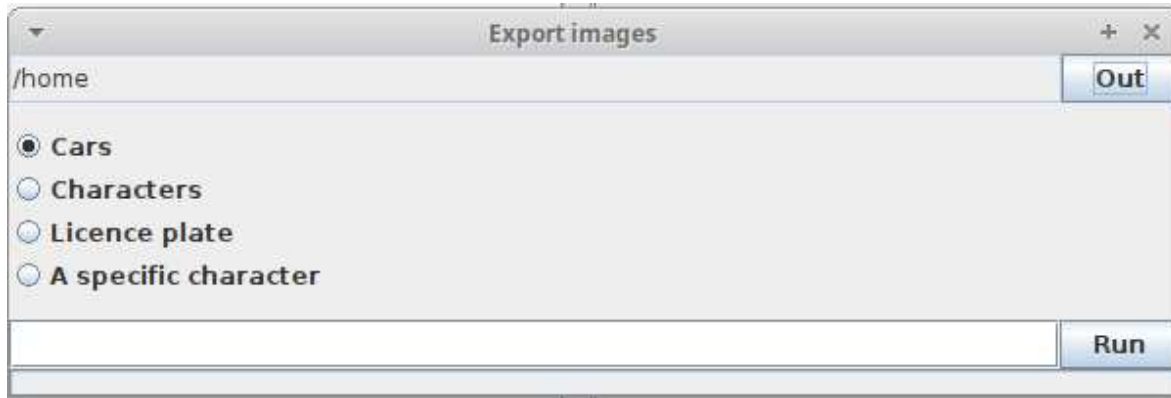


Fig 4. Ventana para exportar las imágenes guardadas

Por último el programa permite observar las estadísticas de las patentes, caracteres y autos recopilados como se muestra en la Fig. 5.

Character	Total quantity	New charactes	Old characters
A	45	39	6
B	15	9	6
C	12	11	1
D	8	6	2
E	10	2	8
F	8	6	2
G	4	0	4
H	8	3	5
I	4	1	3
J	13	4	9
K	7	1	6
L	8	4	4
M	12	3	9
N	7	2	5
O	8	1	7
P	12	1	11
Q	7	2	5
R	9	2	7
S	5	2	3

Cars: 75
Licence Plates: 74 (New: 32 Old: 42)
Characters & Digits: 475 (New: 201 Old: 166)

Fig. 5: En esta ventana se puede observar cuántos caracteres, dígitos, autos y patentes nuevas y viejas fueron recolectadas

La Fig. 5. muestra la cantidad de elementos que fueron segmentados al comenzar con la etapa de reconocimiento. Si bien la cantidad segmentada es magra , fue imperioso comenzar a investigar sobre re-

des neuronales que realmente puedan detectar objetos en tiempo real. De cualquier forma este acervo puede crecer muy rápidamente debido a que ya se encuentran confeccionados todos los programas pertinentes

Como se mencionó anteriormente, teniendo ya una pequeña cantidad de imágenes para poder utilizar nos adentramos en investigar qué redes neuronales son capaces de detectar correctamente en tiempo real elementos en un video. Nuestra primera aproximación es entonces validar de alguna forma las imágenes ya obtenidas entrenando una red neuronal y en caso de funcionar correctamente utilizarla como apoyatura, si fuera posible, para hacer la segmentación manual más rápida.

Como grupo de investigación ya conocíamos las redes neuronales convolucionales como el Neocognitron y las redes LeNet y en parte su principio de funcionamiento, pero desde el año 2015 una red neuronal particular comenzó a ser realmente muy práctica, utilizada y precisa para las tareas de reconocimiento de imágenes en tiempo real. Esta nueva red neuronal, convolucional, se denomina YOLO (You Only Look Once). El principio de funcionamiento es complejo ya que YOLO define una arquitectura de red neuronal convolucional particular que trata de predecir en qué lugar se encuentra el o los objetos en la imagen, para esto YOLO escala la imagen a un determinado ancho y alto (esto anchos y altos están definidos por las distintas resoluciones de redes YOLO, ya que la entrada de esta red neuronal convolucional es toda la imagen) y la divide en segmentos. YOLO busca el segmento que mayor probabilidad posee para determinada clase y devuelve un rectángulo que contiene el objeto reconocido en la imagen. Esta red neuronal fue evolucionando con el paso del tiempo, y en el año 2018 la versión 3 de YOLO fue presentada (YOLOv3) y la versión 4 de YOLO fue presentada en abril del año 2020. En teoría la red neuronal es perfecta para nuestra primera etapa de segmentación de autos y quizás también podría funcionar para detectar las patentes con la imagen del automóvil. Para comprobar que realmente esta red neuronal es realmente funcional, lo primero que realizamos fue intentar hacer funcionar un ejemplo de esta red para conocer su eficiencia.

La red YOLOv3 comúnmente se ofrece pre entrenada con la una base de datos de patrones llamada COCO (Common Objects in Context), esta red pre entrenada reconoce 80 tipos de objetos distintos como personas, perros, motos y entre éstos también automóviles, por lo tanto fue la red neuronal perfecta para probarla. Después de una difícil implementación obtuvimos resultados muy satisfactorios con en una imagen estática de 1920x1080 píxeles utilizando una red YOLO cuya entrada es de 416x416, sin embargo la detección en videos en tiempo real (de la misma resolución que la imagen estática) era demasiado lenta debido a que la inferencia se realizaba solamente en la CPU sin utilizar la GPU en lo absoluto, aunque la eficiencia en la detección era igualmente buena (Fig. 6). Cabe destacar que esta prueba se realizó sobre una computadora con microprocesador I7-7700HQ con 16 GB de memoria y una placa de video Nvidia GTX 1050.

Evidentemente el uso de una GPU disminuiría considerablemente el tiempo de detección de los objetos en YOLO por lo tanto después de instalar todas las dependencias pertinentes pudimos ejecutar en otra prueba esta red neuronal convolucional utilizando la GPU integrada en la computadora y los resultados fueron realmente muy alentadores. mientras que utilizando solo la CPU se alcanzaban valores de entre 5 a 6 fps utilizando la GPU se alcanzaron valores de entre 14 a 15 fps

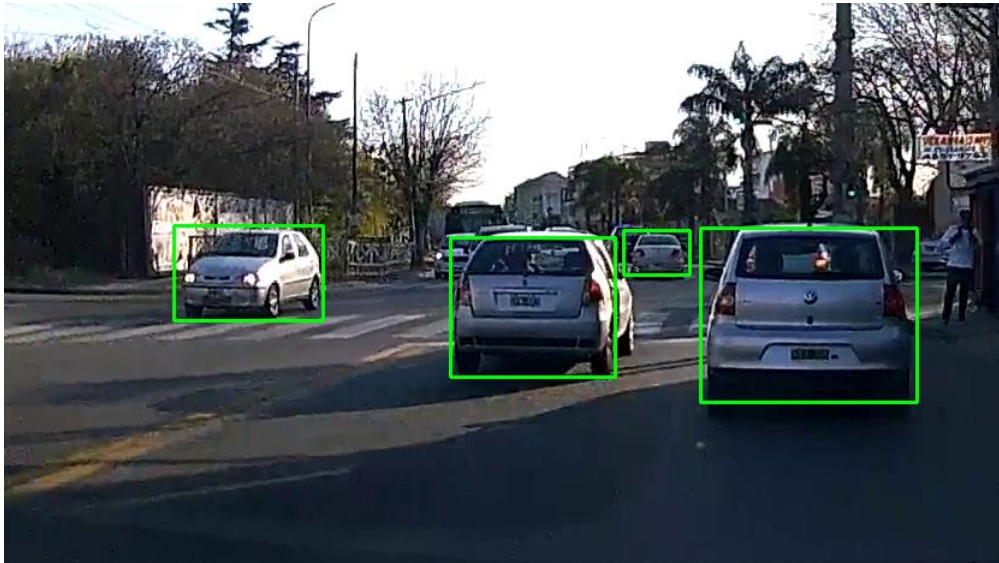


Fig. 6. Fotograma en donde una red neuronal YOLOv3 320 (pre entrenada con la base de datos de COCO) detecta correctamente 4 vehículos

También es importante recalcar que esta no es la red neuronal más precisa en la detección de objetos , pero si es realmente rápida cuando se habla de detección de objetos en tiempo real. Un ejemplo de la falta de precisión de esta red lo pudimos observar en la detección de una camioneta , la red detecta dos automóviles en vez de una sola camioneta. Fig. 7.



Fig 7. Una detección fallida de la red YOLO, a la derecha la ampliación del fallo

Este tipo de errores repercuten demasiado a la detección , ya que la camioneta se encuentra alejada, pero en ciertas ocasiones , el auto negro que se encuentra a la derecha y al frente de la imagen , algunas veces no es detectado correctamente, aunque esto puede ser completamente compensado entrenando mejor la red o agregando un poco más de lógica al detector

Si bien la red YOLOv3 es realmente buena en la detección en tiempo real de objetos parte de nuestro objetivo es implementar todo el sistema de detección de imágenes en tiempo real en un sistema embebido, como por ejemplo una Raspberry 2 / 3 , y si bien la ejecución de YOLOv3 cumple con las expectativas esperadas , la realidad es que la computadora en la que se ejecutó la prueba dista por mucho a las características de hardware que ofrece un sistema embebido como la placa anteriormente mencionada. Sin embargo para esto, ya existe un diseño de una arquitectura de YOLOv3 adaptada para sistemas embebidos que se denomina YOLO-Tiny, que como su nombre lo indica es una versión reducida de la red neuronal convolucional YOLO así que ya con un pequeño acervo de datos , nos dedicamos a entrenar una red YOLO en su versión más pequeña con el conjunto de datos obtenidos de automóviles. Después de 15 horas de entrenamiento de la red YOLO (40000 iteraciones) los resultados no fueron en lo absoluto alentadores, ya que los rectángulos que marcaban los automóviles eran muy grandes respecto a la del objeto encontrado haciendo que la precisión de la detección sea realmente muy deficiente. La figura 8 muestra esta situación. Esto se debe a la forma en la que la red neuronal aprende cuales son los tamaños de los rectángulos que deben contener al objeto al encontrarlo, cuando entrenamos la red neuronal con las imágenes obtenidas del acervo , estas imágenes contenían solamente al objeto en cuestión pero YOLO también necesita la información de la imagen completa y de la relación entre el ancho del objeto que



Fig 8. Una detección fallida después de entrenar la red YOLO por primera vez

que va a ser utilizado para entrenar y el ancho de la imagen donde se encuentra dicho objeto , el alto del objeto que va a ser utilizado para entrenar y el alto de la imagen donde se encuentra dicho objeto, la posición absoluta respecto en la imagen , etc. Al tener ya el acervo creado, en lugar de rearmar el mismo te-

niendo en cuenta los datos necesarios para entrenar la red neuronal YOLO, creamos un pequeño programa para encontrar los recortes de imágenes de los automóviles en los fotogramas de los videos que fueron utilizados para crear el acervo, de esta forma ya obtuvimos las nuevas imágenes y los datos necesarios para poder entrenar la red YOLO. Si bien ahora los rectángulos son acordes con los automóviles , y el entrenamiento es mucho más rápido , después de 20000 iteraciones la red neuronal no tiene una buena performance, aunque es atribuible al magro acervo de datos.

Conclusión:

Debido a los diversos inconvenientes que se presentaron en esta investigación, los cuales fueron motivos ajenos a este grupo, y que golpearon con fuerza a la misma, la investigación no se pudo concluir o finalizar por completo, sin embargo los distintos puntos que se buscaban alcanzar, como la creación de un acervo de datos y la detección de patentes en automóviles no se encuentra muy lejana de llevarse a cabo por completo: El acervo de datos puede crecer rápidamente debido a las herramientas desarrolladas, y si fue posible detectar automóviles con la red YOLO , también será posible entrenar a la red neuronal para detectar patentes en los autos , tanto patentes nuevas como viejas. Probablemente esto implique un desafío de implementación, pero el objetivo probablemente pueda ser alcanzado. Otro punto a investigar es el comportamiento de redes neuronales similares para realizar la tarea de detección de objetos en tiempo real y poder compararla con YOLO , y por último resta investigar si la utilización de la red YOLO es útil en la detección de caracteres en tiempo real , o existe alguna alternativa para un problema que en teoría ya no es una desafío para el campo de la inteligencia artificial debido a las diversas redes neuronales que a lo largo de los años se desarrollaron para cumplir con este objetivo.

Anexo II

Desarrollo de un sistema de asistencia visual para un Manipulador automático programable (robot industrial)

Ing. Martín Ferreyra Biron, Ing. Alberto Miguens, Ing. Fernando Szklanny
Departamento de Ingeniería e Investigaciones Tecnológicas
Universidad Nacional de La Matanza
zamferreyra@unlam.edu.ar, amiguens@unlam.edu.ar, fszklanny@unlam.edu.ar

Resumen

El presente trabajo hace referencia al desarrollo de un sistema de procesamiento digital de imágenes, basado en cámaras de alta velocidad, para incorporar en un autómata programable, capaz de moverse articuladamente en seis ejes, y destinado a aplicaciones industriales con exigencia de eficiencia.

El autómata, que actualmente se encuentra funcionando prácticamente en un 100%, está formado por tres elementos: un manipulador (o "brazo"), un controlador y un conjunto de sensores. El controlador contiene el sistema procesador y las unidades de control de energía. El brazo es accionado por un servomecanismo y es capaz de moverse articuladamente en seis ejes. Los sensores utilizados, con la finalidad de mejorar el control del autómata y la seguridad en el entorno de trabajo, son cámaras de alta velocidad y acelerómetros.

El fundamento del presente trabajo consiste en el desarrollo del sistema de procesamiento digital de imágenes que queda soportado en el autómata en cuestión.

Introducción

La aplicación de los sistemas digitales del procesamiento de señales y su aplicación al control de motores de inducción y autómatas programables tiene su origen en la década de 1970, con los primeros desarrollos de procesadores digitales de señales DSP por Intel, AMI y Bell Labs. Las dificultades en cuanto a la capacidad de procesamiento y velocidad para realizar cálculos, hacen que su aplicación se limite a controles de baja complejidad. No obstante, el gran avance tecnológico producido desde ese momento, tanto en la miniaturización de los sistemas digitales, la capacidad de cálculo en los procesadores de la época y su velocidad de procesamiento hacen que se puedan aplicar a controles con algoritmos de alta complejidad.

El autómata propuesto es un brazo mecánico capaz de moverse articuladamente en seis ejes, del tipo de "articulación coordinada", llamado así por la semejanza

de los movimientos con los del cuerpo humano. La programación de los movimientos del autómata se realiza desde una computadora personal mediante un software a desarrollar.

Los autómatas programables, tal como se mencionó en el apartado anterior, se utilizan ampliamente en la industria automotriz. Pero existen numerosos sistemas de producción que no cuentan con autómatas que se adapten a sus necesidades, por lo que el trabajo que aquí se describe permite proponer soluciones innovadoras que cubran dichas necesidades.

En particular se hace indispensable el desarrollo de algoritmos para el guiado de autómatas en base a la visión, lo que permitirá utilizarlos para detectar y manipular objetos de diferentes tamaños, color y contraste, mejorar la precisión al realizar una trayectoria, y evitar objetos no previstos en una tarea, haciendo al autómata adaptable a diferentes condiciones de trabajo.

En conclusión, con el fin de mejorar la productividad de las industrias, la calidad y el costo de fabricación para ganar competitividad a nivel global se hace interesante contar con autómatas programables en los procesos de fabricación con sistemas de visión.

Marcoteórico

El procesamiento digital de imágenes ha tomado un gran impulso en la última década, derivado de la alta capacidad de los procesadores que posibilitan el desarrollo de sistemas más versátiles y de menor costo.

Por ejemplo, en el control de procesos realizados por robots industriales, la navegación en vehículos autónomos, la detección de eventos e inspección de objetos en procesos de manufacturas y en un sin fin de aplicaciones más.

En la actualidad existe una gran cantidad de bibliografía sobre estudios y algoritmos desarrollados para el control de motores de inducción, sistemas de control electromagnéticos, vehículos eléctricos, procesamiento digital de imágenes y autómatas programables con una vasta aplicación en la industria. Sin embargo, existen numerosos sistemas de producción que no cuentan con

autómatas que se adapten a sus necesidades de proceso automatizados para un aumento del nivel de producción y calidad, menos aunque utilicen el procesamiento digital de imágenes y sea por su elevado costo o por la inexistencia en el mercado. A partir de esta situación se decidió desarrollar un brazo robot que pueda ser automatizado y tomar decisiones en su movimiento, nutriéndose de la entrada que proporcionan las cámaras de video.

Para que el brazo pueda tomar decisiones necesitará reconocer objetos, como se mencionó anteriormente, y esto implica la adecuada elección de un algoritmo que cumpla con este fin. Al momento de iniciar la investigación sobre estos algoritmos se comenzaron a realizar pruebas con SURF [1] y ORB [2] pero no se lograba una correcta detección en objetos que no tuvieran textura. Tampoco se deseaba utilizar redes neuronales convolucionales [3] debido al costo de procesamiento y de entrenamiento más allá de su buen desempeño en reconocer patrones y la atención [4], junto con las redes neuronales recurrentes no se encontraban en ese momento tan difundidas en el campo del procesamiento digital de imágenes.

De todas las opciones estudiadas la que más se acerca a nuestros requerimientos fue el algoritmo DOT [5] que fue el implementado y el cual se detalla en los siguientes párrafos.

Implementación de DOT

Como parte del proyecto desarrollado se planteó el desarrollo e implementación del trabajo de Stefan Hinterstoisser y otros, en el que se presenta un método (DOT) para detectar objetos tridimensionales en tiempo real, bajo una cierta cantidad de ruido, que no requiere una etapa de entrenamiento demasiado costosa y que puede detectar objetos sintéticos (problema que, si se observó, de manera empírica, al utilizar SURF como se mencionó).

DOT

DOT es un algoritmo de detección de objetos presentado en la conferencia de visión por computadora y reconocimiento de patrones del año 2010. Este algoritmo tiene la capacidad de poder detectar los objetos de una manera muy rápida y repetitiva para lo cual se utilizan distintas imágenes del mismo objeto, a las que se denominan patrones (a partir de ahora se utilizará la terminología original "templates") del objeto a buscar. Debido a la forma en que se generan estos "templates" y al método con el que se busca el objeto en la escena, este algoritmo es robusto en la detección de objetos sintéticos. DOT se puede resumir en tres grandes pasos:

- Creación de los "templates".
- Preparación de la escena.
- Búsqueda del objeto.

A continuación, y rápidamente se realizará una breve descripción de cada uno de los pasos anteriormente mencionados

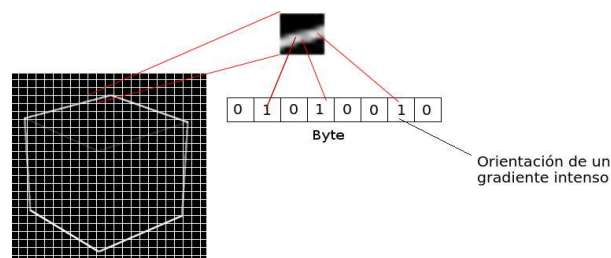


Figura 1: Ejemplo de cómo se realiza un template



Figura 2 : Templates del envase de un adhesivo vinílico

Creación de templates

La primera etapa es la creación de los "templates" a partir de los objetos a buscar. La metodología propuesta por DOT es simple: obtener imágenes del objeto desde distintas perspectivas. Mientras más imágenes distintas se obtengan para generarlos "templates", más eficiente será la búsqueda, pero también más lenta. Estas imágenes serán la materia prima de los "templates" y se buscarán posteriormente en la escena. Un ejemplo de cómo se realiza un template se puede observar en la figura 1.

Posteriormente, a cada "template" se le aplica el filtro Sobel [6] para así poder hallar los gradientes de la imagen. La imagen resultante se divide en cuadrados "pequeños" de $N \times N$ píxeles y se recogerá la orientación de los 7 primeros gradientes que tengan más intensidad en cada cuadrícula. Estos datos se guardarán en un byte siendo el MSB el bit que indique si en esa celda de la cuadrícula hubo o no gradientes. Un ejemplo de templates se puede observar en la figura 2.

Preparación de la escena

Se denomina "escena" a la imagen en la que se buscará el objeto a detectar. El proceso es similar a como se prepara un "template": los bordes de la escena se detectan con un filtro Sobel para así obtener los gradientes de esta, luego se divide la imagen en cuadrados de la misma medida utilizada en los "templates" con la diferencia de que no se buscan los 7 gradientes más intensos por cada uno de los cuadrados, sino que se guarda la orientación del gradiente más intenso. En el caso de trabajar con video esto se tiene que hacer fotograma a fotograma y la manera en que se debe programar el algoritmo tiene que ser muy rápida para ofrecer una detección acorde a los tiempos del brazo robot.

Búsqueda

La etapa de búsqueda es conceptualmente muy sencilla. Se tienen en cuenta las orientaciones de los gradientes más fuertes de cada cuadrado en el “template” y la orientación del gradiente más fuerte de cada cuadrícula de la escena. Según la investigación DOT, se debe deslizar el “template” sobre la escena, calculando cuántas orientaciones del mismo coinciden con las orientaciones de la escena. En el lugar donde hubo mayores coincidencias es donde se supondría que se encuentra el objeto a detectar.

Primera implementación

En esta investigación el algoritmo de DOT fue implementado en más de una ocasión y en diversas formas buscando siempre una optimización en velocidad. La primera implementación de DOT fue realizada en un computador a arquitectura x86 sin tener en cuenta ningún tipo de optimización. Si bien los resultados fueron alentadores, desde el punto de vista del comportamiento

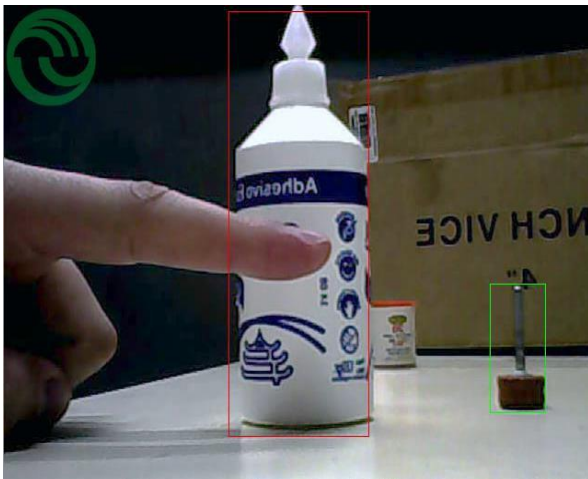


Figura 3: Detección de un envase de adhesivo vinílico con una pequeña occlusión y una piedra de desgaste.

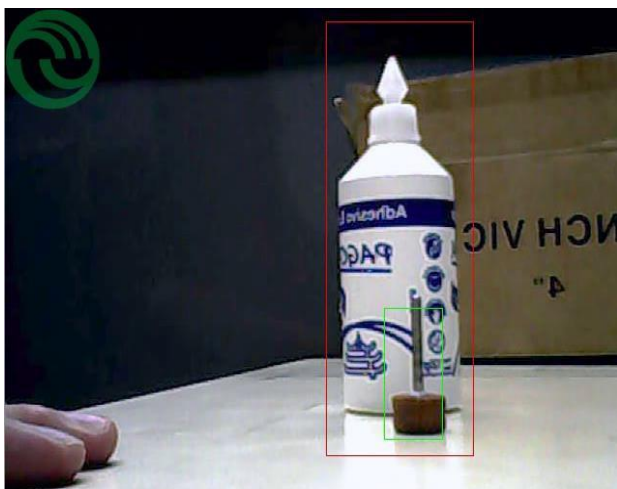


Figura 4: Dos elementos detectados mientras uno oculta al otro

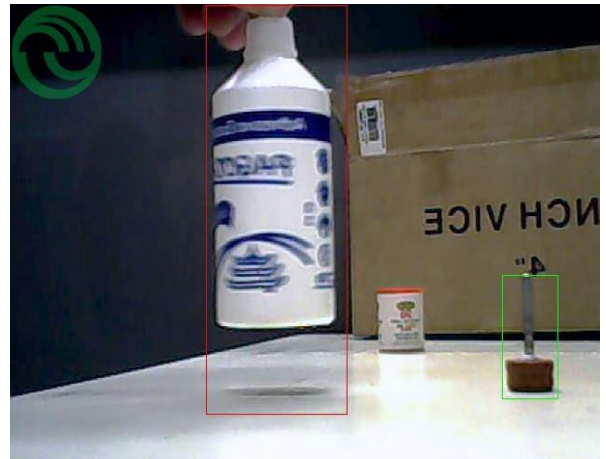


Figura 5: Detección en movimiento

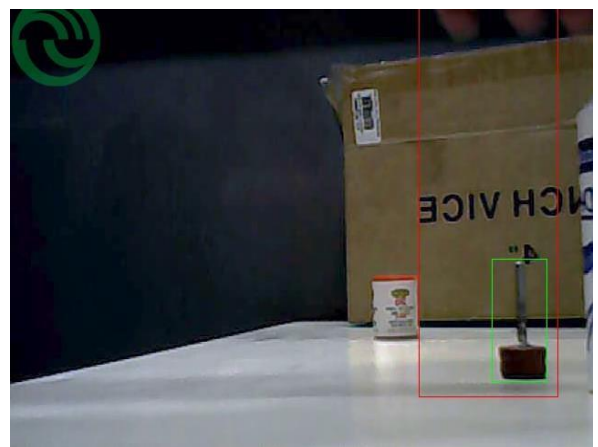


Figura 6: Error en la detección

del algoritmo, fue necesario programarlo nuevamente con mucho más cuidado.

En la segunda implementación, también programada para una arquitectura x86, se tuvo en cuenta la optimización del código y se trató de implementar, como recomiendan los autores de DOT [7], las instrucciones MMX. Con todos estos cambios, se obtuvo un aumento en la performance del algoritmo, pero la misma no fue la suficiente para poder utilizarlo en un brazo robot. Sin embargo, se pudo programar el algoritmo de una forma mucho más eficiente como se detallará en la continuación.

En la figura 3 se puede observar la detección de un envase de adhesivo vinílico y una piedra de desgaste. En el caso de los envases se puede detectar con cierto nivel de occlusión sin mayores inconvenientes. En el caso de la figura 4 se puede observar como el algoritmo puede detectar un objeto ocuyendo al otro, siendo ambos detectados correctamente.

En la figura 5 se puede observar cómo se detecta el envase de adhesivo vinílico en movimiento. En la figura 6 el algoritmo programado falla al no encontrar el objeto. Este es uno de los problemas a resolver.

Implementación de DOT en GPU

Con el objetivo de avanzar en el desarrollo de algoritmos para la detección y procesamiento de imágenes, se procedió a la incorporación al sistema de una placa de desarrollo Nvidia Jetson TK1. La placa en cuestión está construida sobre la base de un procesador ARM@Cortex™-A15 y una unidad destinada al procesamiento gráfico de NVIDIA Kepler GPU con 192 núcleos CUDA. Cabe mencionar que CUDA es una arquitectura de cálculo paralelo de NVIDIA que aprovecha la potencia de la GPU (unidad de procesamiento gráfico) para proporcionar un incremento del rendimiento del sistema.

En un primer intento se compiló el código realizado anteriormente con mínimas modificaciones para ejecutarlo en la placa de desarrollo. Las consideraciones iniciales, basadas en que gran parte del código hace uso de la biblioteca OpenCV, [8] y en que varios métodos utilizados están optimizados para la ejecución en una GPU, prometían o sugerían un buen rendimiento. Pero esto no fue así. La detección en un video de prueba superaba por mucho el tiempo de procesamiento al desuejecución mediante el análisis en una CPU.

Visto y considerando este inconveniente se comenzó a analizar el código con detenimiento. Se observó que la manera en la que se estaba intentando utilizar la GPU no era la correcta, por lo que se examinó cada uno de los métodos utilizados para identificar aquellos que producían la mayor demora y optimizarlos. A partir de este análisis se llegó a la conclusión de que el cálculo de gradientes era el primero en quemarse y demoraba en ejecutarse. Este método se debe ejecutar fotograma a fotograma.

En el cálculo de gradientes las tareas a realizar son las siguientes y cada una de ellas depende de la anterior:

1. Se convierte a escala de grises la imagen del fotograma a analizar.
2. Se calcula el filtro Sobel por X e Y.
3. Se calcula el gradiente a partir de los filtros anteriores.
4. Se halla el máximo gradiente en cuadrados de $N \times N$ píxeles en toda la imagen.
5. Se calcula el ángulo de dicho gradiente guardando su orientación en un byte.
6. Se almacena el resultado.

De toda la lista anteriormente mencionada, el primer lugar en donde se atacó la demora fue en el cálculo del filtro Sobel. Se utilizó la función optimizada para CUDA de OpenCV de manera correcta y debido a que se tienen que hallar los componentes en X y en Y, estas tareas se ejecutaron en paralelo. Sin embargo, si bien se obtuvo una mejora sustancial, la comprensión y el análisis de los tiempos globales de demora llevó a los responsables de la investigación a la conclusión de que la mejor manera de

implementar el algoritmo DOT no era utilizando los méto-

dos optimizados ofrecidos por OpenCV, sino a través del desarrollo de un kernel. Esto fue producto de dos circunstancias, en un principio, desconocidas al momento de iniciar las pruebas:

- El tiempo de transmisión de datos desde la memoria principal a la memoria de la GPU (y viceversa) no es despreciable y de manera inherente genera una sobrecarga en la ejecución del algoritmo ("overhead"). Por lo tanto, utilizar una GPU sintener en cuenta lo anterior produce demoras y las mismas deben ser minimizadas.

- Se debe buscar en la imagen máximos locales en regiones de interés de $N \times N$ píxeles. En primer lugar, OpenCV no ofrece una función optimizada que utilice CUDA para realizar esta tarea y en el caso de ofrecerla se debería tener en cuenta el tiempo de transferencia entre la memoria principal y la memoria de la GPU.

Por todo esto se procedió al diseño de un kernel que pudiera acelerar y hacer uso correcto de la GPU y así hallar los gradientes de manera rápida y correcta. El enfoque utilizado para resolver el problema fue dividir y solucionar el mismo en etapas parciales.

Lo primero que se realizó fue el desarrollo de un kernel que permitiera aplicar el filtro Sobel a una imagen utilizando la GPU. Este paso fue fundamental ya que en un kernel se solucionaban los puntos 2 y 3 de la lista de tareas mencionada anteriormente.

Seguidamente se ensayó el comportamiento del kernel sobre un video y los resultados fueron muy alentadores, en un tiempo mucho menor de lo esperado la aplicación del filtro era calculado. Además, se comprobó también que la placa aceptara el uso de una webcam la cual se utilizó con el algoritmo y funcionó de manera correcta. Por lo tanto, la lista de tareas que hasta ahora se resume es la siguiente:

1. Se convierte a escala de grises la imagen del fotograma a analizar.
2. Se calcula el filtro Sobel por X e Y el gradiente (GPU).
3. Se halla el máximo gradiente en cuadrados de $N \times N$ píxeles en toda la imagen.
4. Se calcula el ángulo de dicho gradiente guardando su orientación en un byte.
5. Se almacena el resultado.

El paso siguiente consistió en la búsqueda del máximo gradiente en cuadrados de $N \times N$ píxeles en la GPU.

La primera aproximación para solucionar este problema fue programar en un kernel distinto la búsqueda de máximos locales utilizando la GPU en bloques de $N \times N$ píxeles. De esta manera se trabajó la imagen resultante del filtro anterior en cuadrados $N \times N$, pudiéndose hallar paralelamente los máximos de cada bloque que. Esto fue correcto pero el hecho de utilizar un único hilo para buscar el máximo gradiente desperdiciaba el poder de procesamiento de la GPU. Es por esto que después de investigar y analizar el problema se reprogramó la solución utilizando el concepto de reducciones paralelas, en este caso, en matrices.

La idea de las reducciones paralelas es sencilla: Para hallar un máximo se deben utilizar comparaciones. Si se utiliza un solo núcleo o un solo hilo, ese único hilo o núcleo debe realizar todas las comparaciones. Pero al tener una GPU y poder decidir cuantos hilos se ejecutan por cada bloque, esta situación cambia. Como la operación de comparación es binaria, la máxima cantidad de hilos que pueden utilizarse para solucionar el problema de manera óptima, es la mitad de la cantidad de elementos (si ésta es par).

Por ejemplo, si se trata de hallar el máximo en una matriz de 8x8 elementos existirán 64 elementos a comparar y si se utilizan 32 hilos se podrán comparar 32 pares de elementos de una sola vez, guardando el resultado en una de las mitades. Si esta operación se continúa realizando (con los elementos que sobran para comparar) la eficiencia que se obtiene es realmente considerable.

La figura 7 ejemplifica de una manera clara el concepto de reducción que se programó.

Posteriormente se halló el máximo, en el mismo kernel se programó el cálculo para transformar y expresar su orientación medida en un byte, tal como lo proponen los autores del algoritmo DOT. Por lo tanto, la lista de tareas que se resume hasta ahora es la siguiente:

1. Se convierte a escala de grises la imagen del fotograma a analizar.
2. Se calcula el filtro Sobel por X e Y, y el gradiente (GPU)
3. Se halla el máximo gradiente en cuadrados de $N \times N$ píxeles en toda la imagen, se calcula el ángulo de dicho gradiente guardando su orientación en un byte y se almacena el resultado (GPU).

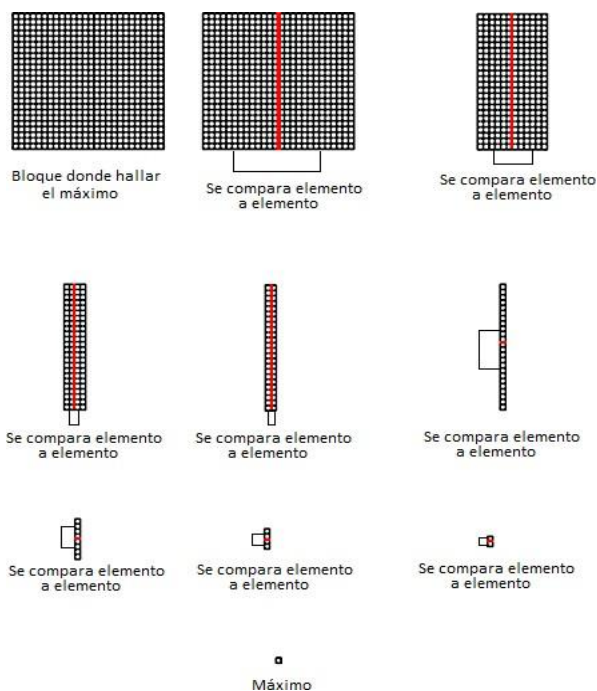


Figura 7: Ejemplificación de la reducción de una matriz utilizando una GPU

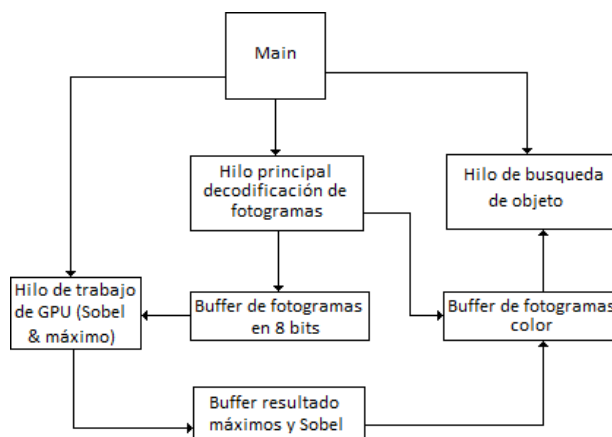


Figura 8: Esquema de ejecución del algoritmo

Por último, debido a que existían demoras entre la ejecución de ambos kernels, estos se combinaron en uno solo. Todas estas mejoras se dejaron a un aumento de velocidad con si de abaratar respecto a las mismas tareas se ejecutadas en una CPU o utilizando en parte la GPU, a través de OpenCV.

Ya solucionado el problema generado por la demora del método que calcula los gradientes, se siguió analizando el programa en búsqueda de otras demoras. Un punto clave que producía las mismas, fue la decodificación del método de compresión utilizado en los videos o en la misma cámara de la que se obtienen los fotogramas. Debido a que la única manera de aumentar la velocidad de esta decodificación es utilizar la GPU, y la transferencia entre la memoria principal y la de la GPU produce demoras, se decidió decodificar los fotogramas en un hilo.

Se debe tener en cuenta que la placa Nvidia Jetson TK1 posee un microcontrolador ARM de cuatro núcleos y utilizar uno en particular para decodificar los fotogramas no produce ningún tipo de problemas.

Valiéndose de una arquitectura de productor y consumidor, en el que en un hilo se decodifican los fotogramas (productor) y en otro hilo se ejecuta el kernel donde se buscan los gradientes y los máximos (consumidor), se aumentó aún más la eficiencia de ejecución del algoritmo DOT. El esquema de la figura 8 muestra cómo se comunican y ejecutan los diversos algoritmos.

Solamente resta optimizar la búsqueda del objeto con la GPU denominado "Match Template".

Optimización de la búsqueda

En primer lugar, se tiene que no se puede sobrecargar más a la GPU en este punto, por lo tanto, la búsqueda de objetos debe realizarse en la CPU. El algoritmo DOT plantea que lo que se debe realizar para encontrar un objeto es buscar la coincidencia más probable entre la escena y los objetos que se aprendieron. Para poder hacer esto se necesita utilizar una ventana deslizante (que recorre toda la escena) y buscar coincidencias de ángulos de los objetos aprendidos. El ganador o la región candidata a ser el objeto en cuestión es aquella en la que existan más coincidencias de orientaciones entre la plantilla del objeto aprendido y la escena. En una versión simplificada del código se vería así:

Por cada plantilla de objeto aprendido
 Por cada pixel en Y de la escena
 Por cada pixel en X de la escena
 Por cada pixel en Y de la plantilla
 Por cada pixel en X de la plantilla
 $Match += Escena(x,y) \& plantilla(x,y)$

Es evidente que la complejidad computacional es alta y de alguna manera habría que minimizarla. Después de analizarla forma en que funcionaba el algoritmo, se llegó a la siguiente conclusión:

Se podría aumentar la eficiencia del algoritmo si solamente se comprobara la coincidencia de ángulos en los lugares donde en la plantilla están definidos, de esa manera se evitaría una iteración "for" y, por ende, varias operaciones en la búsqueda. Además se notó, experimentalmente, que es mucho más rápido hacer la búsqueda cada dos píxeles en la escena que uno por uno (en este caso se estaría bajando la resolución de búsqueda). Para lograr esto se modificó el método que utiliza el programa de aprendizaje para analizar y guardar los objetos aprendidos, ahora en vez de guardar toda la plantilla se guarda solamente el valor del ángulo y la posición en la misma. En resumen, una versión simplificada del código se vería así

Por cada plantilla del objeto aprendido
 Por cada pixel en Y de la escena
 Por cada pixel en X de la escena
 Por cada ángulo válido en la plantilla
 If (Escena(x,y) es un ángulo definido)
 $Match += Escena(x,y) \& plantilla(x,y)$

Si bien el cambio es drástico, es evidente que existen mejoras.

Cambios drásticos en los objetos.

El algoritmo DOT no es perfecto y en las pruebas se observó más de una vez que quizás por falta de plantillas de objetos la detección no se realizaba en el lugar correcto. Para mitigar este comportamiento se puede suavizar de alguna forma la detección brusca de objetos, teniendo en cuenta que los mismos no pueden aparecer en otro lugar en la escena de una forma inmediata, sino que necesitan trasladarse por la misma. Para solucionar este problema se implementó la siguiente solución: Si se detectó un objeto en una posición que supera umbral de detección, ya habiéndose detectado el mismo, la ventana de detección se trasladará de manera progresiva al lugar donde se supone que se encuentra el objeto.

Pruebas en ambiente real

A continuación, se presentan algunas imágenes de momentos en los que el algoritmo produjo resultados apropiados, en videos que se tomaron en ambiente real con ruido de fondo. Es importante recalcar que fueron algunos momentos debido a que no se uti-

lizaron demasiadas plantillas para realizar la detección. Las imágenes de las figuras 9 a 16 son grises debido a que se consumió demasiada memoria y se utilizaron los tres canales de colores. El exceso de consumo de memoria genera, como consecuencia, que el sistema operativo termine de manera abrupta el proceso de detección.



Fig. 9. Algoritmo reconociendo el envase de pegamento



Fig. 10. Algoritmo presentando una falla en la detección



Fig. 11. Algoritmo detectando un cubo correctamente



Fig. 12. Algoritmo fallando por falta de plantillas. En este video en particular las de fallas fueron muy altas. El algoritmo malinterpretaba el ruido de fondo (rompecabezas)



Fig. 13. Tomando un video donde se reconocía el cubo en la implementación en CPU se trató de reconocerlo. En rojo la detección de esta implementación y en azul de la anterior. Cabe destacar que se utilizaron diferentes plantillas

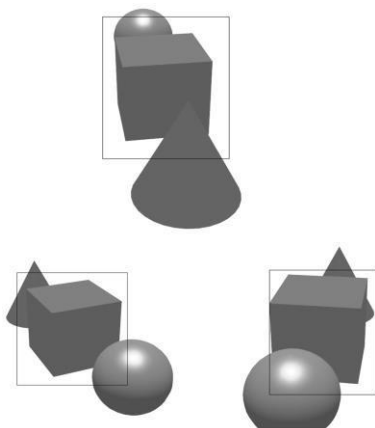


Fig. 14. Detecciones correctas en el video de prueba que se utilizó. Observar la detección con conclusiones

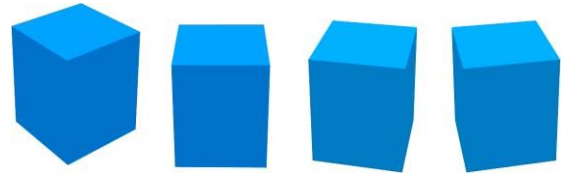


Fig.15. Plantillas que se utilizaron para la detección de la imagen anterior. El video de prueba no tuvo falsos positivos.

Detección de profundidad

La detección de un objeto a partir de la implementación de DOT permite encontrar el mismo en el plano, pero no en el espacio, siendo parte fundamental la detección de esta dimensión debido a la naturaleza de la posición del brazo robot. Es por esto que se necesitó realizar la medición de la profundidad por medio de un sistema estereoscópico. La posibilidad que se evaluó es la de ubicar dos cámaras en un mismo plano M (con coordenadas x y y) separadas por una distancia b en el eje x , y con sus ejes focales paralelos y perpendiculares al plano M . De esta forma se lograrían dos imágenes disparas a nivel horizontal, y el producto de esa disparidad se utilizaría para medir la profundidad.

A nivel matemático [9], este mecanismo puede ser resuelto a través del método de triangulación. En la figura 16 se puede observar un esquema básico, en el cual se muestran dos sistemas de ejes cartesianos cada uno correspondiente a cada cámara del par estereoscópico.

A una profundidad dada por la distancia focal de las cámaras f (que debe ser la misma en ambas) se encuentra el plano denominado epipolar que contiene los planos I_1 e I_2 correspondientes a las imágenes tomadas por cada cámara. En este plano es donde se miden las disparidades entre dichas imágenes.

Al encontrarse las cámaras a la misma altura (eje y), las únicas diferencias se dan a nivel horizontal ya que, como se comentó, las cámaras se ubican a una distancia de separación b sobre el eje x .

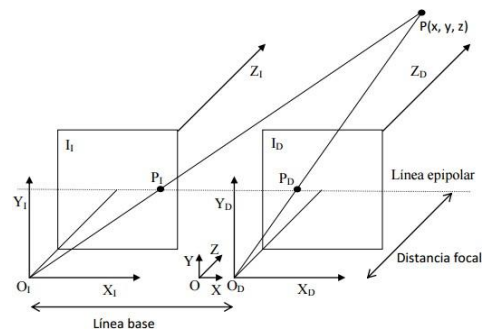


Figura 16 :Diagrama tridimensional para el cálculo de la distancia del punto P

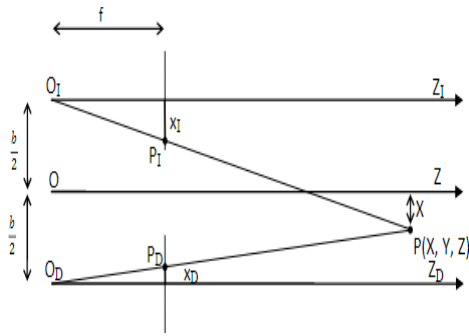


Figura 17: Diagrama bidimensional para el cálculo de la distancia del punto P

En la figura 17, se pueden observar los puntos PI y PD que son las proyecciones del punto P sobre la línea epipolar. Se puede intuir que cuanto más lejano se encuentre el punto P, los vectores de proyección tenderán a ser paralelos al eje Z. Lo que hará PI y PD tiendan a ubicarse sobre el eje ZI y ZD.

La disparidad del punto se calcula como $X_I - X_D = d$ y en el caso de que un objeto se encuentre ubicado en el infinito de eje Z haría que d valga 0. Así entonces el cálculo de la distancia del objeto sobre eje Z será

$$\text{Imagen Izquierda: } \frac{\frac{b}{z} + x}{f} = \frac{x_I}{f} \Rightarrow x_I = \frac{(x + \frac{b}{z})f}{z}$$

$$\text{Imagen Derecha: } \frac{-x}{z} = \frac{x_D}{f} \Rightarrow x_D = \frac{(x - \frac{b}{z})f}{z}$$

En la ecuación *Imagen Derecha* se asigna un signo negativo a X_D para compatibilizar los sistemas de coordenadas de los ejes OI y OD. Siendo la disparidad del punto P: $x_I - x_D = d$ reemplazando y despejando se llega a la ecuación:

$$z = \frac{fb}{d}$$

Con la cual se puede calcular la distancia z del punto P (objeto detectado)

Por otra parte, el desarrollo de los algoritmos de control del automatizado se apoyó en las bibliotecas desarrolladas por Peter Corke, realizándose simulación mediante el software MATLAB.

Queda pendiente la integración con las bibliotecas que se desarrollaron para la detección de objetos, debido a que se priorizó el desarrollo de las bibliotecas de procesamiento gráfico.

Calibración de las cámaras

Lo ideal es que, en un principio, el ajuste se realizara a grandes rasgos lo más precisamente posible, para que luego el algoritmo de calibración tuviera mejores resultados.

Una vez hecho este ajuste, se comenzó con el cálculo de las matrices de calibración, para lo que se utilizó la biblioteca de OpenCV que contiene clases para calibración y reconstrucción de imágenes 3D: "Camera calibration and 3D reconstruction". Para la calibración se usó la aplicación *calibrate_cameras*, que forma parte de la biblioteca Stereo Vision (Bajo licencia pública GNU)

Esta aplicación toma como argumentos imágenes estereó que contienen en su escena un patrón cuadrado como un tablero de ajedrez, con capturas en distintas perspectivas de este (figura 18). El tablero se rotó y además se lo trasladó a lo largo de las múltiples tomas de imágenes. Matemáticamente se puede demostrar que la cantidad mínima de tomas depende de la cantidad de esquinas que contiene el tablero. El término "esquinas" hace referencia a la cantidad de vértices de cuadrados negros que hacen contacto entre sí (léase negro y blanco). La máxima cantidad utilizada es 50, bajo el criterio de que cuanto más muestras haya los cálculos de corrección de errores serán más precisos.

La aplicación *Calibrate_cameras* da como resultado un archivo con un conjunto de matrices que contienen los parámetros intrínsecos (focos, centrados), extrínsecos (traslación y rotación) y las matrices de rectificación de distorsiones. Con estas matrices se realizan las correcciones sobre las imágenes tomadas en cada cámara. Una vez rectificadas las imágenes se puede realizar el cálculo de eje Z con la fórmula citada más arriba.

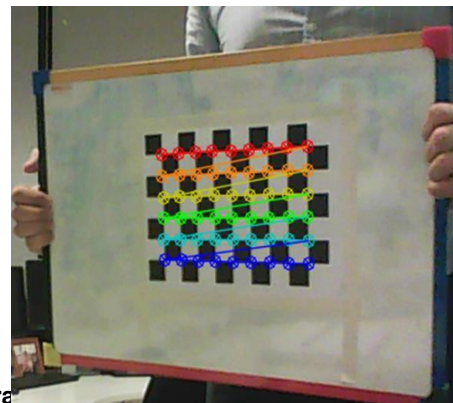


Figura 18: Ejemplo de un tablero de ajedrez utilizado para la calibración "findChessboardCorners" de la biblioteca OpenCV

Para lograr este cálculo hay que ubicar en ambas imágenes del par estéreo el mismo objeto, para esto es necesario un algoritmo de búsqueda e identificación de objetos. Esto se puede resolver con el algoritmo DOT que fue detallado anteriormente.

Verificación con mapa de disparidad

La verificación con mapa de disparidad se realizó utilizando las bibliotecas de *StereoVision*, una prueba empírica realizando un mapa de disparidad y una nube de puntos. La nube de puntos consiste en conseguir las coordenadas espaciales de cada pixel de la escena tomada por las cámaras. De esta forma se puede visualizar una escena desde diferentes perspectivas rotando los ejes de coordenadas. Esto tiene un límite y es que esos "puntos" de la imagen se obtienen desde una posición fija de espacio, si se desea observar perspectivas laterales se necesitaría más información de la escena. Esto se puede salvar hasta cierto ángulo, replicando los píxeles de la frontera, de forma que cubran los huecos de información faltante.

Para lograr un buen mapa de disparidad se necesita que los objetos de la escena sean detectados inequívocamente en las imágenes del par estéreo. Se utilizó la herramienta *tune_blockmatcher* de la biblioteca *StereoVision* con el fin de ajustar los parámetros del detector de objetos y se verificó en paralelo la "calidad" del mapa de disparidad que se muestra en la figura 19.

Los ajustes mencionados se deben realizar en forma iterativa y empírica ya que resulta ser el método más rápido y eficiente. Luego de realizar el ajuste del algoritmo de búsqueda y detección de objetos, se utilizó la función *images_to_pointcloud* (también de la citada biblioteca *StereoVision*), la que toma como argumentos las matrices de rectificación, un par de imágenes estereo, los parámetros de ajuste del algoritmo de "blockmatching" y esto da como resultado una matriz de tres dimensiones que contiene los vectores (también de tres dimensiones) correspondientes a la información del color de cada pixel. Como resultado se puede observar la siguiente nube de puntos (Figuras 20 y 21)

Conclusión

Quizás las modificaciones y los agregados realizados sobre el algoritmo DOT no sean suficientes para obtener velocidades considerables en la detección, aunque permiten mejorar adecuadamente la misma. Ahora es aceptable realizar una búsqueda en una imagen HD, pero

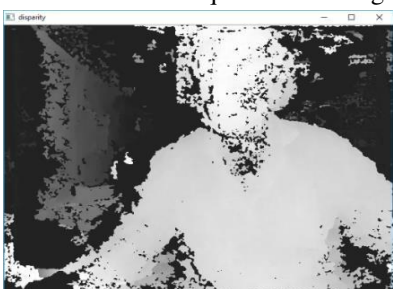


Figura 19: Mapa de disparidad



Figura 20: Un par de imágenes (imagen estereo)

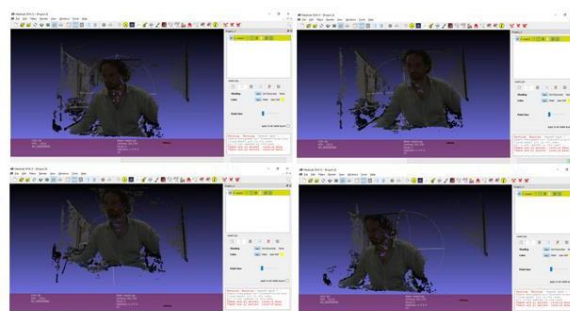


Figura 21: resultado de la nube de puntos, a partir de la imagen estereo.

aun así sigue siendo relativamente lento. Como se plantea al final de la parte "Implementación DOT en GPU" quizás la respuesta involucra la implementación del proceso de "match template" en la GPU, aprovechando el tiempo ocioso de la misma.

Respecto al algoritmo, presenta fallas que quizás sean mejorables. DOT no considera colores, solo considera formas, y queda claro que resolver esa falencia sería una mejora a esta implementación. Otro punto importante a destacar es que el tiempo de reconocimiento de un objeto depende de la cantidad de plantillas y de la dimensión del video. Esto no debería ser así, aunque es una característica intrínseca del algoritmo.

Una solución que se podría proponer para la mejora del algoritmo es la de no utilizar la ventana deslizante, sino que por algún medio se detectaran los posibles candidatos a la imagen de la plantilla buscada y luego realizar el "match template" entre las plantillas y esa región de la imagen.

Respecto al cálculo y análisis de profundidad, la calibración de las cámaras y la comprensión del algoritmo están en un nivel muy avanzado, al punto de que solamente sería necesario programar la detección y agregarla a la de objetos. El único inconveniente será el tiempo de procesamiento que se sumaría al que posee la detección DOT y, como se expresó, se debería buscar una solución de compromiso o analizar la unión de ambos algoritmos mucho más detenidamente.

En etapas siguientes de este proyecto de investigación, que sigue vigente en el ámbito de la Universidad, se procederá a continuar el trabajo sobre aquellas cuestiones que hoy se mencionan como dificultades o asuntos pendientes.

Referencias

- [1] Bay, Herbert, Tuytelaars, Tinne, Van Gool, Luc.: SURF: Speeded Up Robust Features, European Conference on Computer Vision 2006
- [2] Rublee, Rabaud, Konolige, Bradski: ORB: an efficient alternative to SIFT or SURF, IEEE International Conference on Computer Vision 2011
- [3] Lecun, Haffner, Bottou, Bengio, Object Recognition with Gradient-Based Learning, CPVR 2004
- [4] Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, Polosukhin, Attention is all you need, Cornell University 2017
- [5] Hinterstoisser Stefan, Lepetit Vincent, Ilic Slobodan, Fua Pascal, Navab Nassir.: Dominant Orientation Templates for Real-Time Detection of Texture-Less Objects, Computer Vision Pattern Recognition 2010.
- [6] Gonzalez, Woods: Digital Image Processing, Tercera Edición, Prentice Hall 2007
- [7] Hinterstoisser Stefan, Cagniard Cedric, Ilic Slobodan, Sturm Peter, Navab Nassir, Fua Pascal, Lepetit Vincent.: Gradient Response Maps for Real-Time Detection of Texture-Less Objects, IEEE Transactions on Pattern Analysis and Machine Intelligence 2011
- [8] OpenCV <http://opencv.org/>
- [9] Martín Montalvo, Técnicas de visión estereoscópica para determinar la estructura tridimensional de la escena, Autor: Curso académico 2009/2010, Máster en Investigación en Informática, Facultad de Informática, Universidad Complutense de Madrid 2009/2010.

Artículo original

DESARROLLO DE UN SISTEMA DE ASIS- TENCIA VISUAL PARA UN MANIPULADOR AUTOMÁTICO PROGRAMABLE (ROBOT INDUSTRIAL)

VISUAL ASSISTANCE SYSTEM FOR A PROGRAMMABLE AUTOMATIC HANDLER (INDUSTRIAL ROBOT)

Ing. Martín Ferreyra Biron, Ing. Alberto Miguens, Ing. Fernando Szklanny

Departamento de Ingeniería e Investigaciones Tecnológicas, Universidad Nacional de La Matanza
mferreyra@unlam.edu.ar , amiguens@unlam.edu.ar , fszklanny@unlam.edu.ar

DESEA PUBLICAR SU E-MAIL (tachar lo que no corresponda): (NO)

Resumen:

En el presente trabajo se describe el desarrollo de un sistema de procesamiento digital de imágenes y la adaptación de un algoritmo detector de objetos, utilizando cámaras de alta velocidad con el objetivo final de incorporar dicho sistema a un autómatas programable, apto para moverse con seis grados de libertad de forma articulada y con el fin de ser utilizado en aplicaciones industriales con ciertas exigencias de eficiencia. El autómatas, que se encuentra funcionando casi en su totalidad, está conformado por tres elementos: un manipulador, un controlador y un conjunto de sensores.

El controlador es el encargado de procesar la información para mover el manipulador y las unidades de control de energía. Por otra parte, el brazo está comandado por un servomecanismo que permite mover al mismo articuladamente en seis ejes, y por último los sensores utilizados, con la finalidad de mejorar el control del autómatas y la seguridad en el entorno de trabajo, son cámaras de alta velocidad y acelerómetros.

El objetivo principal del presente consiste en el desarrollo del sistema de visión que dará soporte al autómatas descrito.

Abstract:

This document describes the development of a digital image processing system and an object detector adaptation using high-speed cameras, with the final goal of being included into a programmable automaton, suitable to move in six degrees of freedom, and oriented to industrial applications with specific efficiency requirements. This robot, currently working almost entirely, includes three elements: a manipulating system (the arm), a controller and a quantity of sensors.

The controller is the device that processes the movement information and the energy control units. The arm is moved by a servomechanism and is capable of moving in six axes and, finally, the sensors required, in order to get a better control of the robot and to ensure safety in the work environment, are high-speed cameras and accelerometers.

In other terms, the main objective of the current work consists in the development of a visual system that will support the abovementioned robot.

Palabras Clave: *Visión, robot industrial, velocidad, seguridad.*

Key Words: *Vision, Industrial robot, speed, safety.*

I. INTRODUCCIÓN

La aplicación de los sistemas digitales de procesamiento de señales y su aplicación al control de motores de inducción y autómatas programables tiene su origen en la década de 1970, con los primeros desarrollos de procesadores digitales de señales DSP por Intel, AMI y Bell Labs. Las dificultades en cuanto a la capacidad de procesamiento y velocidad hicieron que su aplicación se limite a controles de baja complejidad. Sin embargo, el avance tecnológico producido desde ese momento hasta hoy, hace factible aplicar estos sistemas digitales de procesamiento a controles con algoritmos de alta complejidad.

El autómata propuesto es un brazo mecánico capaz de moverse articuladamente en seis ejes, del tipo de "articulación coordinada", llamado así por la semejanza de sus movimientos con los del cuerpo humano. La programación de los movimientos del autómata se realiza desde una computadora personal mediante un software a desarrollar. Los autómatas programables se utilizan ampliamente en la industria automotriz, aunque existen numerosos sistemas de producción que no cuentan con autómatas que se adapten a sus necesidades, por lo que el trabajo que aquí se describe permite proponer soluciones innovadoras que cubran dichas necesidades. En particular se hace indispensable el desarrollo de algoritmos para el guiado de autómatas en base a la visión, lo que permitirá utilizarlos para detectar y manipular objetos de diversas índoles, mejorar la precisión al realizar una trayectoria, y evitar objetos no previs-

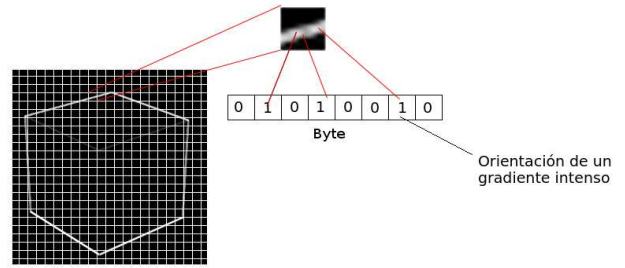


Fig.1. Ejemplo de cómo se realiza un "template"

tos en una tarea, haciendo al autómata adaptable a diferentes condiciones de trabajo.

En conclusión, con el fin de mejorar la productividad de las industrias, la calidad y el costo de

fabricación para ganar competitividad a nivel global, se hace interesante contar con autómatas programables en los procesos de fabricación con sistemas de visión.

II. MARCO TEÓRICO

El procesamiento digital de imágenes ha tomado un gran impulso en la última década y si bien en la actualidad existe una gran cantidad de bibliografía sobre estudios y algoritmos desarrollados para el control de motores de inducción, sistemas de control electromagnéticos, vehículos eléctricos, procesamiento digital de imágenes, entre otros, no abundan sistemas de producción que cuenten con autómatas que se adapten a sus necesidades respecto a procesos automatizado-orientado a un aumento del nivel de producción y calidad y menos aún que utilicen el procesamiento digital de imágenes. A partir de esta situación se decidió desarrollar un brazo robot que pueda ser automatizado y tomar decisiones en sus movi-

mientos, valiéndose de la entrada que proporcionan cámaras de video. Para que el brazo pueda tomar decisiones necesitará reconocer objetos y esto implica la adecuada elección de un algoritmo que cumpla con tal fin. Al inicio de la investigación sobre estos algoritmos se hizo un estudio empírico sobre SURF[1] y ORB[2] los cuales forman parte de conocidos detectores aunque no se lograba una correcta detección en objetos que no tuvieran textura. Tampoco se deseaba utilizar redes neuronales convolucionales [3] debido al costo de procesamiento y de entrenamiento más allá de su buen desempeño en reconocer patrones y la atención [4], junto con las redes neuronales recurrentes no se encontraban en ese momento tan difundidas en el campo del procesamiento digital de imágenes, por lo tanto de todas las opciones estudiadas la que más se acercaba a nuestros requerimientos fue el algoritmo DOT [5] que fue el implementado y el cual se detalla en los siguientes párrafos.

III. IMPLEMENTACIÓN DE DOT

Como se comentó en párrafos anteriores, parte del proyecto implica el desarrollo e implementación de DOT, algoritmo del cual daremos una breve reseña.

A. DOT

DOT es un algoritmo de detección de objetos presentado en la CVPR 2010 el cual tiene la capacidad de poder detectar objetos con o sin textura de una manera rápida y repetitiva valiéndose de distintas imágenes procesadas del objeto a buscar, a las que se denomina patrones (a partir de ahora se utilizará la terminología original “*templates*”). DOT se puede resumir en tres grandes pasos:

- Creación de los “*templates*”.
- Preparación de la escena.
- Búsqueda del objeto.

B. CREACIÓN DE TEMPLATES

La primera etapa es la creación de los “*templates*” a partir de los objetos a buscar, siendo la metodología propuesta por DOT simple: obtener imágenes del objeto desde distintas perspectivas para luego procesarlas y obtener finalmente los “*templates*”. Mientras más “*templates*” distintos se obtengan, más eficiente será la búsqueda, pero también más lenta. Estas imágenes serán la materia prima de los “*templates*” y se buscarán posteriormente en la escena. Luego, a cada imagen del objeto obtenida se le aplica el filtro Sobel [6] para así poder hallar los gradientes de la imagen. La imagen resultante se divide en “pequeños” cuadrados de $N \times N$ pixel y se recogerá la orientación de los primeros 7 gradientes más intensos en cada cuadrícula. Estos datos se guardarán en un byte siendo el MSB el bit que indique si en esa celda de la cuadrícula hubo o no gradientes.

C. PREPARACIÓN DE ESCENA

Se denomina “escena” a la imagen en la que se buscará el objeto a detectar. El proceso es similar a como se prepara un “*template*”: los bordes de la escena son detectados con el filtro Sobel para así obtener los gradientes, luego se divide la imagen en cuadrados de la misma medida utilizada en los “*templates*” con la diferencia de que no se buscarán los 7 gradientes más intensos por cada cuadrícula, sino que se reservará únicamente la orientación del gradiente más intenso.

D. BÚSQUEDA

Finalmente, la etapa de búsqueda, conceptualmente sencilla, consta en “deslizar” sobre la escena cada “template” calculando cuántas orientaciones de los gradientes más intensos de la escena coinciden con las orientaciones de los “templates”. En el lugar donde haya mayores coincidencias es donde se supondría que se encuentra el objeto a detectar.

IV. IMPLEMENTACIÓN DE DOT EN GPU

Si bien la implementación propuesta en el trabajo de Hinterstoisser fue realizada en una computadora con arquitectura x86 sin GPU, la performance obtenida no fue la esperada para poder aplicar esta técnica a un sistema de visión de un brazo robot, por lo tanto, con el objetivo de avanzar en el desarrollo del algoritmo, se procedió a la incorporación al sistema de una placa de desarrollo Nvidia Jetson TK1 la cual posee una GPU.

En un primer intento se compiló el código realizado anteriormente con mínimas modificaciones para adaptarlo a la placa mencionada. Las consideraciones iniciales, sumado a que gran parte del código hace uso de la biblioteca OpenCV, [8] y que varios métodos utilizados están optimizados para la ejecución en una GPU, sugerían un buen rendimiento, pero esto no fue así. La detección en un video de prueba superaba por mucho el tiempo de procesamiento al de su ejecución mediante el análisis en una CPU. Observando este inconveniente se comenzó a analizar el código con detenimiento. Se observó que la manera en la que se estaba intentando utilizar la GPU no era la adecuada, por lo tanto, se examinó cada uno de los métodos utilizados para identificar aquellos que producían la mayor demora y optimizarlos. A par-

tir de este análisis se llegó a la conclusión de que el cálculo de gradientes en la escena era el primero que más demoraba en ejecutarse. Recordemos que el método se debe ejecutar fotograma a fotograma ya que se trata de un video. En el cálculo de gradientes las tareas a realizar son las siguientes y cada una de ellas depende de la anterior:

1. Se convierte a escala de grises la imagen del fotograma a analizar.
2. Se calcula el filtro Sobel por X e Y.
3. Se calcula el gradiente a partir de los filtros anteriores.
4. Se halla el máximo gradiente en cuadrados de $N \times N$ píxeles en toda la imagen.
5. Se calcula el ángulo de dicho gradiente guardando su orientación en un byte.
6. Se almacena el resultado.

De toda la lista anteriormente mencionada, el primer lugar en donde se atacó la demora fue en el cálculo del filtro Sobel. Se intentó utilizar, nuevamente, pero de otra forma, el filtro Sobel optimizado para CUDA ofrecido en OpenCV y ejecutado en paralelo (para obtener los componentes X e Y) pero los tiempos globales de detección nos hizo pensar que la mejor forma de ganar rendimiento al implementar el algoritmo DOT no era utilizando los métodos optimizados ofrecidos por OpenCV, sino a través del desarrollo de un kernel. Esto fue producto de dos circunstancias, en un principio desconocidas, al momento de iniciar las pruebas:

- El tiempo de transmisión de datos desde la memoria principal a la memoria de la GPU (y viceversa) no es despreciable y de manera inherente genera

una sobrecarga de tiempo en la ejecución del algoritmo. Estos tiempos deben ser minimizados.

- Se debe buscar en la imagen máximos locales en regiones de interés de $N \times N$ píxeles ya que OpenCV no ofrece una función optimizada que utilice CUDA para realizar esta tarea.

Por todo lo expuesto se procedió al diseño de un kernel que pudiera acelerar y hacer uso correcto de la GPU. y así hallar los gradientes eficientemente. El enfoque utilizado para resolver el problema fue dividir y solucionar el mismo en etapas parciales.

El primer paso fue el desarrollo de un kernel que permitiera aplicar el filtro Sobel a una imagen utilizando la GPU. Esto fue fundamental ya que en un kernel se solucionaban los puntos 2 y 3 de la lista de tareas mencionada. A continuación, se ensayó el comportamiento del kernel sobre un video y los resultados fueron alentadores: en un tiempo mucho menor de lo esperado, el filtro era aplicado. Por lo tanto, con estos cambios la lista de tareas quedaría de la siguiente forma:

1. Se convierte a escala de grises la imagen del fotograma a analizar.
2. Se calculan el filtro Sobel por X e Y y el gradiente (GPU).
3. Se halla el máximo gradiente en cuadrados de $N \times N$ píxeles en toda la imagen.
4. Se calcula el ángulo de dicho gradiente guardando su orientación en un byte
5. Se almacena el resultado.

El paso siguiente consistió en la búsqueda del máximo gradiente en los bloques de $N \times N$ píxeles

en la GPU. La primera aproximación para solucionar este problema fue programar en un kernel distinto la búsqueda de máximos locales utilizando la GPU en bloques de $N \times N$ píxeles. De esta forma se trabajó la imagen resultante del filtro anterior en cuadrados $N \times N$, pudiéndose hallar paralelamente los máximos de cada bloque. Esto fue correcto, pero el hecho de utilizar

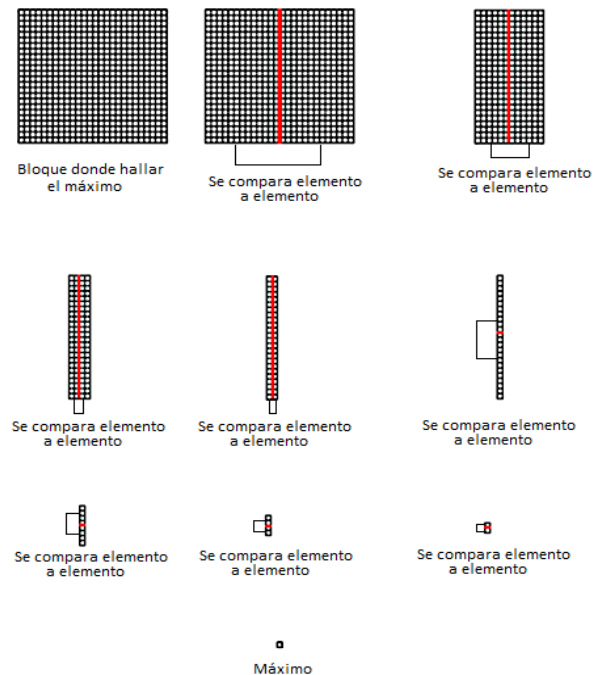


Fig.2. Ejemplificación de la reducción de una matriz utilizando una GPU

un único hilo para buscar el máximo gradiente desperdiciaba el poder de procesamiento de la GPU. Es por esto por lo que después de analizar el problema se reprogramó la solución utilizando el concepto de reducciones paralelas, en este caso, en matrices.

La idea de las reducciones paralelas es sencilla (la figura 2 ejemplifica de una manera clara el concepto de reducción que se programó): para hallar un máximo se deben utilizar comparaciones. Si se utiliza un solo núcleo o hilo, ese único hilo o núcleo debe realizar todas las comparaciones. Pero

al tener una GPU y poder decidir cuantos hilos se ejecutan por cada bloque, esta situación cambia: como la operación de comparación es binaria, la máxima cantidad de hilos que pueden utilizarse para solucionar el problema de manera óptima es la mitad de la cantidad de elementos (si esta es par). Por ejemplo, si se trata de hallar el máximo en una matriz de 8x8 elementos existirán 64 elementos a comparar y si se utilizan 32 hilos se podrán comparar 32 pares de elementos de una sola vez, guardando el resultado en una de las mitades. Si esta operación se continúa realizando (con los elementos que sobran comparar) la eficiencia que se obtiene es realmente considerable.

A posteriori de haberse hallado el máximo, en el mismo kernel se programó el cálculo para transformar y expresar su orientación medida en un byte, tal como lo proponen los autores del algoritmo DOT. Por lo tanto, la lista de tareas con los cambios realizados se resume en:

1. Se convierte a escala de grises la imagen del fotograma a analizar.
2. Se calculan el filtro Sobel por X e Y, y el gradiente (GPU)
3. Se halla el máximo gradiente en cuadrados de NxN pixeles en toda la imagen, se calcula el ángulo de dicho gradiente guardando su orientación en un byte y se almacena el resultado (GPU)

Las optimizaciones descritas producen un considerable aumento de velocidad respecto a las mismas tareas ejecutadas en una CPU o utilizando en parte la GPU, a través de OpenCV.

Ya solucionado en gran parte la demora producida por el cálculo de los gradientes, se continuó

analizando el programa en búsqueda de otras demoras. Un punto clave que producía las mismas, fue la decodificación del método de compresión utilizado en los videos o en la misma cámara de la que se obtienen los fotogramas. Debido a que la única manera de aumentar la velocidad de esta decodificación es utilizar la GPU, y la transferencia entre la memoria principal y la de la GPU produce demoras, se decidió decodificar los fotogramas en un hilo común en la CPU.

Se debe tener en cuenta que la placa Nvidia Jetson TK1 posee un microcontrolador ARM de cuatro núcleos y utilizar uno en particular para decodificar los fotogramas no produce ningún tipo de problemas. Valiéndose de una arquitectura de productor consumidor, en el que en un hilo se decodifican los fotogramas (productor) y en otro hilo se ejecuta el kernel donde se buscan los gradientes y los máximos (consumidor), se aumentó aún más la eficiencia de ejecución del algoritmo DOT.

A. OPTIMIZACIÓN DE LA BÚSQUEDA

En este puntose entiende que la GPU ya no puede ser sobrecargada, por lo tanto, la búsqueda del objeto se debe realizar en la CPU. El algoritmo DOT propone que para encontrar un objeto es necesario buscar la coincidencia más probable entre la escena y los objetos que se aprendieron. Para lograr esto se necesita utilizar una ventana deslizante (que recorre toda la escena) y buscar coincidencias de ángulos de los objetos aprendidos. El “ganador” o la región candidata a ser el objeto en cuestión es aquella en la que existan más coincidencias de orientaciones entre la plantilla del objeto aprendido y la escena. En una versión simplificada del código se vería así:

Por cada plantilla de objeto aprendido

Por cada pixel en Y de la escena

Por cada pixel en X de la escena

Por cada pixel en Y de la plantilla

Por cada pixel en X de la plantilla

Match+=Escena(x,y)&plantilla(x,y)

La alta complejidad computacional salta a la vista y de alguna manera habría que minimizarla. Luego de analizar el algoritmo se llegó a las siguientes conclusiones: se podría aumentar la eficiencia si solamente se comprobara la coincidencia de ángulos en los lugares donde estén definidos en la plantilla, de esa manera se evitaría una iteración “for” y, por ende, varias operaciones en la búsqueda. Además, se notó, experimentalmente, que es mucho más rápido hacer la búsqueda cada dos pixeles en la escena que uno por uno (sacrificando resolución en la búsqueda por velocidad).

Para lograr esto se modificó el método que utiliza el programa de aprendizaje para analizar y guardar los objetos aprendidos, ahora en vez de guardar toda la plantilla se guarda solamente el valor del/los ángulos y la posición en la misma. En resumen, una versión simplificada del código se vería así

Por cada plantilla del objeto aprendido

Por cada pixel par en Y de la escena

Por cada pixel par en X de la escena

Por cada ángulo válido en la plantilla

If(Escena(x,y) es un ángulo definido)

Match+=Escena(x,y)&plantilla(x,y)

Si bien el cambio no es drástico, es evidente que existen mejoras.

B. CAMBIOS DRÁSTICOS EN EL OBJETO

El algoritmo DOT no es perfecto y en las pruebas se observó que quizás por falta de plantillas la detección no se realizaba en el lugar correcto. Para morigerar este comportamiento se intentó suavizar la detección brusca de objetos, teniendo en cuenta que estos no pueden aparecer en otro lugar en la escena de instantáneamente. Esta situación se trató de solucionar de la siguiente manera: si se detectó un objeto en una posición que supera un umbral de detección, la ventana de detección se trasladará de manera progresiva al lugar donde se supone que se encuentra el objeto.

C. PRUEBAS

En la figura 3 se presentan algunas imágenes

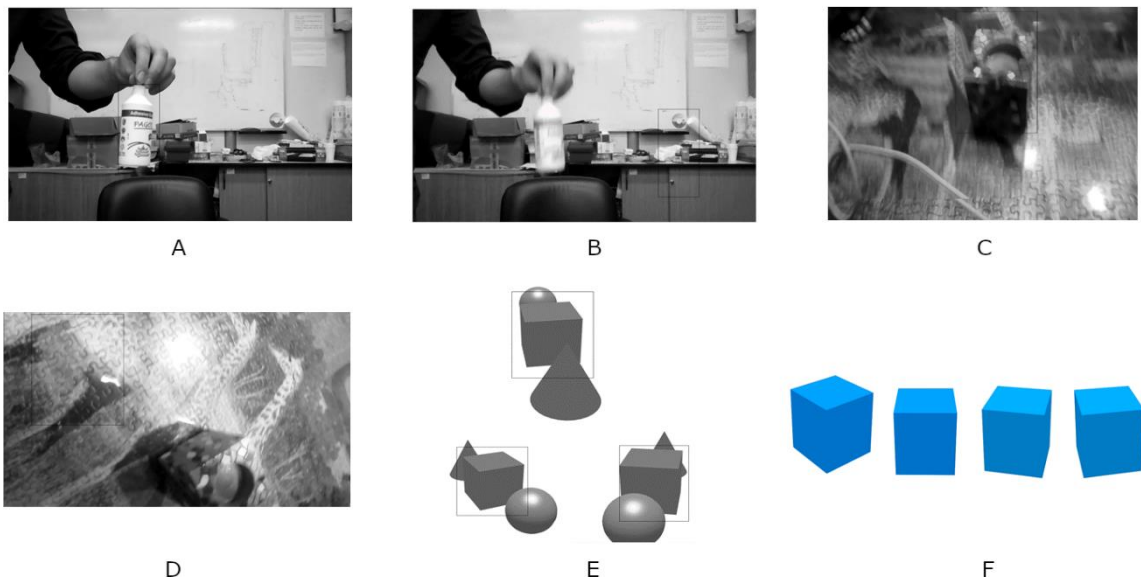


Fig.3. Imágenes resultantes del detector. A) Algoritmo reconociendo el envase de pegamento. B) Algoritmo presentando una falla en la detección C) Algoritmo detectando un cubo correctamente D) Algoritmo fallando por falta de plantillas. E) Detecciones correctas en el video de prueba que se utilizó. Observar la detección con oclusiones. F) Plantillas que se utilizaron para la detección de la imagen E.

nes de momentos en los que el algoritmo produjo resultados apropiados. Es importante recalcar que fueron algunos momentos, debido a que no se utilizaron demasiadas plantillas para realizar la detección. Las imágenes son grises debido a que se consume demasiada memoria si se utilizan los tres canales de colores. El exceso de consumo de memoria genera, como consecuencia, que el sistema operativo termine de manera abrupta el proceso de detección.

V. DETECCIÓN DE PROFUNDIDAD

La detección de un objeto a partir de la implementación de DOT permite encontrar el mismo en el plano, pero no en el espacio, siendo parte fundamental la detección de esta dimensión debido a la naturaleza posicional del brazo robot. Es por esto por lo que se necesita realizar la medición de la profundidad por medio de un sistema estereoscópico. La posibilidad que se evaluó es la de ubicar dos cámaras en un mismo plano M (con coordenadas x e y) separadas por una distancia b en el eje x , y con sus ejes focales paralelos y perpendiculares al plano M . De esta forma se lograrían dos imágenes dispares a nivel horizontal, y el producto de esa disparidad se utilizaría para medir la profundidad.

A nivel matemático [9], este mecanismo puede ser resuelto a través del método de triangulación. En la Figura 4 se puede observar un esquema básico, en el cual se muestran dos sistemas de ejes cartesianos cada uno correspondiente a cada cámara del par estereoscópico.

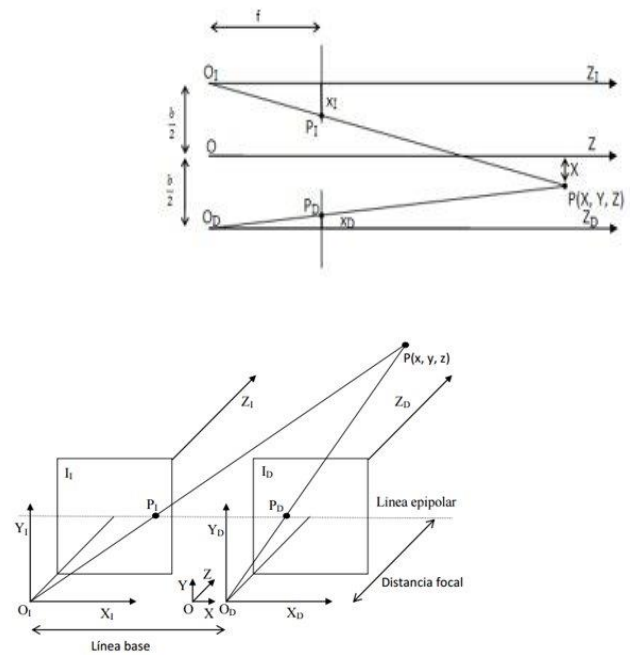


Fig.4. Diagrama tridimensional para el cálculo de la distancia del punto P

A una profundidad dada por la distancia focal de las cámaras f (que debe ser la misma en ambas) se encuentra el plano denominado epipolar que contiene los planos I_1 e I_D correspondientes a las imágenes tomadas por cada cámara. En este plano es donde se miden las disparidades entre dichas imágenes.

Al encontrarse las cámaras a la misma altura (eje y), las únicas diferencias se dan a nivel horizontal ya que, como se comentó, las cámaras se ubican a una distancia de separación b sobre el eje x .

En la Figura 5, se pueden observar los puntos P_1 y P_D que son las proyecciones del punto P sobre la línea epipolar. Se puede intuir que cuanto más lejano se encuentre el punto P , los vectores de proyección tenderán a ser paralelos al eje Z lo que hará que P_1 y P_D tiendan a ubicarse sobre el eje Z_1 y Z_D . La disparidad del punto se calcula como $X_1 - X_D = d$ y en el caso de que un objeto se encuentre ubicado en el infinito del eje Z hará que d valga 0. Así entonces el cálculo de la distancia del objeto sobre el eje Z será:

$$\text{ImagenIzquierda: } \frac{\frac{b}{2} + x}{z} = \frac{x_I}{f} \Rightarrow x_I = \frac{\left(x + \frac{b}{2}\right) f}{z}$$

$$\text{ImagenDerecha: } \frac{\frac{b}{2} - x}{z} = \frac{x_D}{f} \Rightarrow x_D = \frac{\left(x - \frac{b}{2}\right) f}{z}$$

Fig. 5. Diagrama bidimensional para el cálculo de la distancia del punto P

En la ecuación *ImagenDerecha* se asigna un signo negativo a X_D para compatibilizar los sistemas de coordenadas de los ejes O_I y O_D . Siendo la disparidad del punto P: $x_I - x_D = d$ reemplazando y despejando se llega a la ecuación:

$$z = \frac{fb}{d}$$

con la cual se puede calcular la distancia z del punto P (objeto detectado).

Por otra parte, el desarrollo de los algoritmos de control del autómatas se apoyó en las bibliotecas desarrolladas por Peter Corke, realizándose simulación mediante el software MATLAB. Queda pendiente la integración con las bibliotecas que se desarrollaron para la detección de objetos, debido a que se priorizó el desarrollo de las de procesamiento gráfico.

A. CALIBRACIÓN DE LAS CÁMARAS

Lo aconsejable fue, en un principio, que el ajuste se realizara en forma mecánica logrando la mejor alineación posible de las cámaras (ejes horizontal y vertical) para que luego el algoritmo de calibración tuviera mejores resultados. Una vez realizado este ajuste, se comenzó con el cálculo de las matrices de calibración, para lo que se utilizó la biblioteca de OpenCV que contiene clases para calibración y reconstrucción de imágenes 3D: "Camera calibration and 3D reconstruction". Para la calibración se usó la aplicación *calibra-*

te_cameras, que forma parte de la biblioteca StereoVision[10] (Bajo licencia pública GNU).

Esta aplicación toma como argumentos imágenes estereó que contienen en su escena un patrón cuadrículado similar a un tablero de ajedrez, con capturas en distintas perspectivas de éste (Figura



Fig.6. Vértices detectados por el algoritmo de calibración

6). El tablero se rota y, además, se lo traslada a lo largo de las múltiples tomas de imágenes. Matemáticamente se puede demostrar que la cantidad mínima de tomas depende de la cantidad de esquinas que contiene el tablero patrón. El término "esquinas" hace referencia a la cantidad de vértices de cuadrados negros que hacen contacto entre sí (léase negros o blancos). La máxima cantidad utilizada es 50, bajo el criterio de que cuantas más muestras haya los cálculos de corrección de errores serán más precisos.

La aplicación *Calibrate_cameras* da como resultado una carpeta con un conjunto de matrices que contiene los parámetros intrínsecos (focos, centrados), extrínsecos (traslación y rotación) y las matrices de rectificación de distorsiones. Con estas matrices se realizan las correcciones sobre las imágenes tomadas en cada cámara. Una vez rectificadas las imágenes se puede realizar el cálculo del eje Z con la fórmula citada más arriba. Para lograr este cálculo hay que ubicar en ambas imágenes del par estereó el mismo objeto, para esto es necesario un algoritmo de búsqueda e identifica-

ción de objetos. Esto se puede resolver con el algoritmo DOT que fue detallado anteriormente.

B. VERIFICACIÓN CON MAPA DE DISPARIDAD

La verificación con mapa de disparidad se realizó utilizando las bibliotecas de StereoVision, una prueba empírica realizando un mapa de disparidad y una nube de puntos. La nube de puntos consiste en conseguir las coordenadas espaciales de cada pixel de la escena tomadas por las cámaras.

De esta forma se puede visualizar una escena desde diferentes perspectivas rotando los ejes de coordenadas. Esto tiene un límite y es que esos “puntos” de la imagen se obtienen desde una posición fija del espacio, si se deseara observar perspectivas laterales se necesitaría más información de la escena. Esto se puede salvar hasta cierto ángulo, replicando los pixeles de la frontera, de forma que cubran los huecos de información faltante.

Para lograr un buen mapa de disparidad se necesita que los objetos de la escena sean detectados inequívocamente en las imágenes del par estéreo. Se utilizó la herramienta *tune_blockmatcher* de la biblioteca *StereoVision* con el fin de ajustar los parámetros del detector de objetos y se verificó en paralelo la “calidad” del mapa de disparidad que se muestra en la Figura 7.

Los ajustes mencionados se deben realizar en forma iterativa y empírica ya que resulta ser el método más rápido y eficiente. Luego de realizar

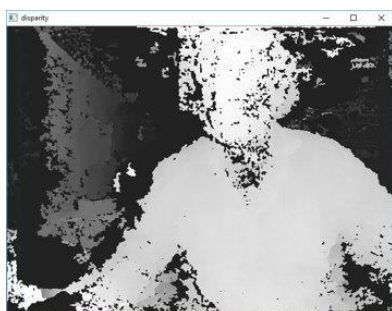


Fig7. Mapa de disparidad

el ajuste del algoritmo de búsqueda y detección de objetos, se utilizó la función *images_to_pointcloud* (también de la citada biblioteca StereoVision), la que toma como argumentos las matrices de rectificación, un par de imágenes estéreo, los parámetros de ajuste del algoritmo de “*blockmatching*”, esto da como resultado una matriz de tres dimensiones que contiene los vectores (también de tres dimensiones) correspondientes a la información del color de cada pixel. Se puede observar la nube de



Fig. 8. Resultado de la nube de puntos a partir de la imagen estéreo

puntos en la Figura 8)

V. CONCLUSIONES

Quizás las modificaciones y los agregados realizados sobre el algoritmo DOT no sean suficientes para obtener velocidades considerables en la detección, aunque permiten mejorar adecuadamente la misma. Ahora es aceptable realizar una búsqueda en una imagen HD, pero aun así sigue siendo relativamente lento. Como se plantea al final del apartado “Implementación de DOT en GPU” quizás la respuesta involucre la implementación del proceso de “match template” en la GPU, aprovechando el tiempo ocioso de la misma.

Respecto al algoritmo, presenta fallas que quizás sean mejorables. DOT no considera colores, solo considera formas, y queda claro que resolver esa falencia sería una mejora a esta implementación. Otro punto importante para destacar es que el tiempo de reconocimiento de un objeto depende de la cantidad de plantillas y de las dimensiones del video. Esto no debería ser así, aunque es una característica intrínseca del algoritmo. Una solución que se podría proponer para la mejora del algoritmo es la de no utilizar la ventana deslizante, sino que por algún medio se detectarían los posibles candidatos a la imagen de la plantilla buscada y luego realizar el “match template” entre las plantillas y esa región de la imagen.

Respecto al cálculo y análisis de profundidad, la calibración de las cámaras y la comprensión del algoritmo están en un nivel muy avanzado, al punto de que solamente sería necesario programar la detección y agregarla a la de objetos. El único inconveniente será el tiempo de procesamiento que se sumaría al que posee la detección DOT y, como se expresó, se debería buscar una solución de compromiso o analizar la unión de ambos algoritmos mucho más detenidamente.

En etapas siguientes de este proyecto de investigación, que sigue vigente en el ámbito de la Universidad, se procederá a continuar el trabajo sobre aquellas cuestiones que hoy se mencionan como dificultades o asuntos pendientes.

VII. REFERENCIAS Y BIBLIOGRAFÍA

A. Referencias bibliográficas:

[1] Bay, Herbert, Tuytelaars, Tinne, Van Gool, Luc.: SURF: Speeded Up Robust Features, European Conference on Computer Vision 2006

[2] Rublee, Rabaud, Konolige, Bradski: ORB: an efficient alternative to SIFT or SURF, IEEE International Conference on Computer Vision 2011

[3] LeCun, Haffner, Bottou, Bengio, Object Recognition with Gradient-Based Learning, CPVR 2004

[4] Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, Polosukhin, Attention is all you need, Cornell University 2017

[5] Hinterstoisser Stefan, Lepetit Vincent, Ilic Slobodan, Fua Pascal, Navab Nassir.: Dominant Orientation Templates for Real-Time Detection of Texture-Less Objects, Computer Vision Pattern Recognition 2010.

[6] Gonzales, Woods: Digital Image Processing , Tercera Edición, Prentice Hall 2007

[7] Hinterstoisser Stefan, Cagniard Cedric, Ilic Slobodan, Sturm Peter, Navab Nassir, Fua Pascal, Lepetit Vincent.: Gradient Response Maps for Real-Time Detection of Texture-Less Objects, IEEE Transactions on pattern analysis and machine intelligence 2011

[8] OpenCV 9/2/2005 [en línea] Fecha de consulta: 27/7/2020 <http://opencv.org/>

[9] Martin Montalvo, Técnicas de visión estereoscópica para determinar la estructura tridimensional de la escena, Autor: Curso académico 2009/2010, Máster en Investigación en Informática, Facultad de Informática, Universidad Complutense de Madrid 2009/2010.

[10] StereoVision 15/04/2017 [en línea] Fecha de consulta 27/7/2020

<https://pypi.org/project/StereoVision/>

Recibido:A completar por el Editor. Formato: AAAA-MM-DD

Aprobado:A completar por el Editor. Formato: AAAA-MM-DD

Hipervínculo Permanente: A completar por el Editor

Datos de edición: Vol. [A completar por el Editor]-Nro. [A completar por el Editor]-Art. [A completar por el Editor]

Fecha de edición: Formato: AAAA-MM-DD



CONMISI

VII Congreso Nacional de Ingeniería
Informática - Sistemas de Información

2019

San Justo, 5 de diciembre de 2019

Se certifica que **Martin Ferreyra Biron, Alberto Raul Miguens, Fernando Ignacio Szklanny** han participado como Autores del artículo 342 "*Desarrollo de un sistema de asistencia visual para un manipulador automático programable (robot industrial)*", aceptado en el VII Congreso Nacional de Ingeniería Informática – Sistemas de Información, CONMISI 2019, realizado los días 14 y 15 de noviembre en la Universidad Nacional de La Matanza.



Ing. Claudio D'Amico
Coord. Gral. CONMISI



Dr. Carlos Neil
Coordinador RII SIC



Mg. Jorge Ezequiel
Decano DIIIT



Ingrese el criterio por el que va a buscar la revista:

(1) ISSN / e-ISSN:

(2) Nombre:

(1) Si conoce el ISSN / e-ISSN de la revista, ingreselo y presione "Buscar por ISSN / e-ISSN" para seleccionarla.
(2) Ingrese una parte del nombre de la revista y presione "Buscar por Nombre" para seleccionarla.

Revista seleccionada:

Editorial:

País de edición: **Ciudad de la editorial:**

Título del artículo:

Idioma:

Volumen: **Tomo:** **Número:**

Página inicial: **Página final:**

(3) Fecha de publicación:

(4) URL:

(4) DOI:

Referato: Sin referato Con referato

(3) Estado de publicación: Publicado En prensa

: