

1 A Generic Framework for Damage-Free Grasping of
2 Delicate Produce Using LLMs and Volume Estimation

3 Ziye Zhang^{b,1}, Xiaoyu Xia^{b,1}, Yuhao Jin^{a,b}, Qizhong Gao^{a,b}, Lin Qiao^{a,b},
4 Jinglei Chen^b, Yong Yue^b, ShanLiang Yao^c, Xiaohui Zhu^{a,*}

5 ^a*Department of Computer Science, University of Liverpool, Liverpool L69 3BX, UK*

6 ^b*School of Advanced Technology, Xi'an Jiaotong-Liverpool University, Suzhou 215123,
7 China*

8 ^c*School of Information Engineering, YanCheng Institute of Technology, Yancheng 224051,
9 China*

10 **Abstract**

Despite notable progress in agricultural robotics, achieving stable, adaptive, and damage-free grasping of irregular, fragile, and diverse produce remains a key unsolved challenge. Traditional approaches rely on geometric heuristics or fixed-force strategies, which cannot accurately model or adapt to the complex physical properties of various produce, often causing instability or damage. To the best of our knowledge, this paper proposes the first adaptive damage-free grasping framework that integrates large language models (LLMs) with multimodal perception to produce. We develop an intelligent framework combining semantic understanding from visual images, geometric awareness from 3D point clouds, and LLMs commonsense knowledge to infer the minimal stable grasping force, enhancing safety and adaptability. To enhance accuracy and robustness, we introduce a point-cloud-based volume estimation module that directly leverages spatial geometry, thereby minimizing reliance on LLMs reasoning and significantly improving perception quality and estimation precision. Moreover, we design a hybrid multi-agent collaboration mechanism that efficiently coordinates perception, reasoning, and control agents, improving reasoning and execution in complex environments. Experiments on a custom dataset of 25 typical produce categories show our framework significantly outperforms state-of-the-art baselines in both volume estimation accuracy and damage-free grasp success rate, demonstrating excellent stability, robustness, and truly damage-free capability. This work validates the effectiveness of combining LLMs and multimodal perception in agricultural robotics and suggests a promising direction for future research toward intelligent and reliable grasping. Our framework introduction is at <https://wyattzzz.github.io/DamageFreeGrasp>

11 **Keywords:** Multimodal Perception, Large Language Model, Agricultural
12 Robot, Damage-free Grasping, Commonsense Physical Reasoning

13 **1. Introduction**

14 Autonomous robotic grasping is critical for agricultural tasks such as produce
15 harvesting and pick-and-place operations (Zhang et al., 2020; Jin and Han,
16 2024). In recent years, deep learning-based grasp detection models, such as
17 GGCNN (Ghasemzadeh et al., 2023), GR-ConvNet (Kumra et al., 2020), and
18 VMGNet (Jin et al., 2024), can efficiently generate candidate grasping poses
19 from visual or point cloud data, providing essential support for stable and robust
20 autonomous grasping (Jin et al., 2025). However, existing grasping strategies
21 still lack sufficient force adaptability, hindering the robot’s ability to dynamically
22 adjust the applied force during the grasping process (Li et al., 2024). First,
23 existing methods lack a generalizable physical reasoning mechanism, making it
24 difficult to adapt to delicate produce with diverse shapes and textures (Zhang
25 et al., 2020; Huang et al., 2025). As a result, the grasping force cannot be
26 precisely controlled, often leading to grasp failure or produce damage (Zheng
27 et al., 2021).

28 To address such challenges, researchers have begun incorporating multimodal
29 information, such as visual, tactile, and semantic data, to enhance the robust-
30 ness and intelligence of grasping strategies (Li et al., 2023b, 2025). Among
31 these, large language models (LLMs) have gained increasing attention for their
32 potential in robotic grasping tasks (Rashid et al., 2023). LLMs can perform
33 commonsense physical reasoning to infer key physical properties of produce, in-
34 cluding volume, density, mass, friction coefficient with the gripper, and elastic
35 modulus (Xie et al., 2024). This capability enables force-controlled grippers to
36 adjust grasping force more precisely. Consequently, the integration of LLMs
37 holds promise for equipping robots with greater adaptability and generalization
38 when interacting with previously unseen produce (Xu et al., 2024). However,
39 despite the physical reasoning capabilities of LLMs, relying solely on LLMs may
40 lead to limited generalization and reduced accuracy in physical inference, par-
41 ticularly when dealing with non-standard produce shapes (Chu et al., 2025; Xie
42 et al., 2024). To mitigate this limitation, point cloud data is incorporated to
43 assist in volume estimation and correct inference bias, thereby reducing depen-
44 dence on LLMs and improving the accuracy of grasp parameter prediction. This
45 also enhances the framework’s robustness and generalization to diverse produce
46 morphologies.

47 In real-world agricultural environments, which are often complex and task-
48 diverse, a single model may struggle to complete the full loop from perception
49 to decision-making (Jin et al., 2025; Hu et al., 2025a). Therefore, this work
50 further introduces a multi-agent collaboration mechanism that systematically
51 integrates semantic reasoning and geometric estimation modules into a cohe-
52 sive workflow designed for damage-free grasping. This mechanism significantly
53 enhances the efficiency of multi-source information fusion and grasp parame-

*Corresponding author.

Email Address: Xiaohui.Zhu@xjtu.edu.cn

¹These authors contributed equally to this work.

54 ter decision-making in complex environments, offering an effective solution for
55 achieving stable grasping of diverse agricultural produce.

56 Our contributions are summarized as follows:

- 57 (1) To the best of our knowledge, this study is the first to propose a grasp
58 reasoning framework that integrates visual, semantic, and point cloud in-
59 formation to intelligently infer the minimal stable grasping force required
60 for delicate produce, significantly improving safety and adaptability in soft
61 grasping.
- 62 (2) By introducing point cloud-based volume estimation of produce, the frame-
63 work reduces reliance on LLMs commonsense reasoning ability, thereby
64 improving the accuracy of grasp parameter estimation and enhancing the
65 framework's robustness and generalization across diverse produce morpholo-
66 gies.
- 67 (3) A multi-agent collaboration system is incorporated into the damage-free
68 grasping of delicate produce, effectively improving the coordination effi-
69 ciency and general applicability of the framework in complex perception
70 and parameter reasoning tasks.

71 The structure of this paper is organized as follows: In the [Section 2](#), we
72 review the literature on the application of LLMs in robotic grasping, as well as
73 grasping techniques for fragile produce. The [Section 3](#) provides a detailed ex-
74 planation of our proposed framework. The [Section 4](#) presents our experimental
75 setup and results. Finally, [Section 5](#) summarizes the paper and discusses the
76 limitations of our work along with potential future directions.

77 2. Related Work

78 Recent research has explored a variety of approaches to achieving intelligent,
79 stable, and damage-free grasping of delicate produce, propelled by advances
80 in LLMs and multimodal perception technologies. These developments have
81 opened up new opportunities for integrating high-level semantic reasoning with
82 low-level force control, improving the adaptability and safety of soft grasping
83 systems. This section reviews related work from two perspectives: LLMs in
84 robotic grasping and delicate grasping techniques.

85 2.1. Large Language Models in Grasping

86 In recent years, the application of LLMs in robotic grasping and manipu-
87 lation tasks has garnered increasing attention. Benefiting from pretraining on
88 large-scale corpora, LLMs not only possess strong capabilities in natural lan-
89 guage understanding and generation but also exhibit a degree of commonsense
90 physical reasoning ([Wang et al., 2023](#)). This enables them to infer implicit prop-
91 erties of produce, including volume, mass, and grasp stability, based on semantic
92 labels. Building on this idea, RT-Grasp ([Xu et al., 2024](#)) proposes a “Reason-
93 ing Tuning” approach, which incorporates a structured reasoning stage during

94 training, allowing multimodal LLMs to generate numerical predictions, such as
95 grasp pose, grounded in semantic context. This significantly enhances control
96 accuracy in complex tasks. Similarly, SayCan (Ahn et al., 2022) combines lan-
97 guage comprehension with physical constraint reasoning to guide robots through
98 step-by-step execution of intricate tasks. However, these methods often exhibit
99 instability when dealing with irregularly shaped produce with complex internal
100 structures. This is primarily due to the lack of common semantic labels and
101 insufficient training data related to such produce, leading to unreliable physical
102 inference by LLMs.

103 To address these limitations, recent research has explored architectural in-
104 novations to enhance the synergy between multimodal perception and language
105 understanding, which are critical for high-level reasoning in robotic systems.
106 For example, Point-BERT (Yu et al., 2022) and PointCLIP (Zhang et al., 2022)
107 employ Transformer-based architectures to encode point cloud data and align
108 it with textual modalities, improving cross-modal semantic consistency. ULIP
109 (Xue et al., 2023) further extends this idea by learning unified representations
110 of images, point clouds, and text for open-vocabulary 3D understanding. Sim-
111ilarly, LLMs such as GPT-4o and QWen-Plus adopt a vision encoder coupled
112 with a language model, enabling them to jointly reason over visual and textual
113 modalities. These models exhibit remarkable generative and reasoning abili-
114 ties in complex visual-textual scenarios, offering valuable insights for robotic
115 decision-making. Collectively, these advancements highlight the potential of in-
116 tegrating semantic reasoning and geometric perception to enable more intelligent
117 and context-aware decision-making in robotic grasping systems.

118 Existing approaches typically use LLMs to directly estimate produce vol-
119 ume or mass based on semantic priors and commonsense knowledge. However,
120 these methods predominantly infer from semantic labels and are limited in their
121 ability to perceive the precise geometric structure of individual instances. This
122 often results in high estimation errors or unstable predictions, particularly when
123 dealing with irregular produce that exhibits diverse and complex shapes.

124 2.2. Grasping Delicate Produce

125 In agricultural robotic grasping, the irregular shapes, fragile skins, and di-
126 verse varieties of produce present significant challenges, motivating extensive
127 research on perception, strategies, and force control. In the domain of visual
128 perception, a method (Chen et al., 2024) integrates 3D object detection with
129 feedforward control, enabling precise grasping of thin-skinned produce. This
130 effectively reduces damage during the grasping process. To further improve the
131 stability of grasp decision-making, a three-finger soft gripper with integrated
132 force feedback and slip detection was designed for apple harvesting (Chen et al.,
133 2022), which significantly reduces bruising and slippage in field experiments.
134 Considering the difficulties of detecting and grasping small produce, as well as
135 their susceptibility to damage, another study (Visentin et al., 2023) proposed a
136 soft gripper that combines RGB visual feedback and tactile sensing. This design
137 not only improved the detection accuracy and grasp success rate for small pro-

138 duce but also minimized potential damage as much as possible. More recently,
139 a soft gripper equipped with tactile sensing and slip detection was developed
140 ([Liu et al., 2024](#)), demonstrating strong adaptability to delicate produce such as
141 grapes and tomatoes in experimental settings. Despite these advances, current
142 approaches often rely on limited visual or geometric cues, making it difficult
143 to assess the underlying physical properties of produce, which are critical for
144 achieving stable, damage-free grasping.

145 To address this challenge, researchers have developed methods that explicitly
146 incorporate texture and firmness perception to assist grasp force modulation and
147 enable selective, damage-free harvesting. For instance, an integrated vision-
148 tactile sensor was developed to estimate produce texture and firmness during
149 the picking process, thereby assisting in grasp force modulation ([Ma et al.,](#)
150 [2024](#)). Furthermore, a soft gripper combined with vision-tactile modalities was
151 proposed to implement a non-invasive firmness estimation mechanism ([Lin et al.,](#)
152 [2023](#)), which was validated to be effective for handling a wide range of produce
153 categories. In addition, a compact tandem-actuated gripper, which is capable of
154 adaptively switching between suction and compliant mechanical grasping modes
155 based on environmental context, enables selective and damage-free harvesting
156 even under partially occluded conditions ([Velasquez et al., 2024](#)). These studies
157 underscore the critical role of integrating texture and firmness perception with
158 adaptive control strategies to achieve more intelligent, reliable, and context-
159 aware damage-free harvesting in diverse agricultural environments.

160 Existing approaches in agricultural robotic grasping often rely on complex
161 soft gripper designs that integrate multiple sensing modalities. While these de-
162 signs enhance adaptability and reduce produce damage, they frequently require
163 intricate mechanical structures and tight sensor integration, which increases
164 system complexity, fabrication cost, and maintenance demands. On the percep-
165 tion side, traditional models for grasp planning and control are typically large
166 and resource-intensive, relying on heavyweight neural architectures that limit
167 real-time deployment on embedded or mobile robotic platforms.

168 3. Methodology

169 This section outlines the proposed framework for damage-free grasping of
170 delicate produce. It is structured into four subsections, starting with a high-
171 level overview of the framework and followed by detailed explanations of each
172 module.

173 3.1. Framework Overview

174 This study introduces an innovative framework that combines 3D point cloud
175 geometric perception with LLMs for the task of damage-free grasping of deli-
176 cate produce. By combining visual perception and language-based reasoning,
177 the framework significantly enhances adaptability and accuracy in the gras-
178 ping process, as illustrated in [Fig. 1](#). The overall workflow consists of three

179 main modules: (1) a 3D perception module, (2) a physical property reasoning
180 module, and (3) a grasp strategy generation and execution module. In the
181 3D perception module, the system performs geometric shape perception using
182 RGB images and depth maps, extracting produce 3D geometric information to
183 support subsequent grasping decisions. The physical property reasoning mod-
184 ule then combines semantic labels with geometric information, leveraging LLMs
185 to infer physical properties such as density and friction coefficients, providing
186 accurate parameters for force-controlled grasping. Finally, the grasp strategy
187 generation and execution module formulates the grasp strategy based on the
188 reasoning results and executes it with a force-controlled gripper to ensure sta-
189 ble, damage-free grasping. The seamless integration of these modules enables
190 the system to perform stable grasping for delicate produce while effectively pre-
191 venting damage, marking the first attempt to apply the combination of 3D
192 geometric perception and LLMs in agricultural robotics.

193 In the 3D perception module, the system receives RGB images and depth
194 maps of the produce. Combined with a semantic label set including apples,
195 bananas, grapes, and others, multimodal feature fusion and semantic-guided
196 querying are performed using the multimodal fusion encoder and Language-
197 guided Query Selection (LGQS) modules. The fused features are then processed
198 by the Multimodal Decoder, which generates 2D coordinates of the target pro-
199 duce, enabling precise localization and semantic label alignment. Next, the
200 identified produce mask, along with its depth map, is processed through Planar
201 Surface Extraction and Non-Planar Segmentation, extracting 3D contour
202 information and forming a complete point cloud and shape representation of
203 the produce. A coarse-grained volume estimation of produce is then obtained,
204 which will be input to the grasp execution module.

205 In the physical property reasoning module, the system utilizes the semantic
206 label extracted from the 3D perception module, along with information about
207 the gripper material, to drive three LLMs-powered agents called FricAgent,
208 AdjustAgent, and DensAgent, which respectively infer the friction coefficient,
209 volume correction factor, and density estimate of the produce.

210 In the grasping strategy generation and execution module, the module de-
211 termines the minimum stable grasping force using the Contact Force Estimator,
212 based on the coarse volume and inferred physical attributes of the produce,
213 thereby ensuring safe manipulation without damage. Simultaneously, the 2D
214 coordinates and depth map of the produce are processed by the Grasp Pose
215 Generator, an extension of VMGNet (Jin et al., 2024), to predict candidate
216 grasp poses. These poses are jointly optimized with the estimated minimum
217 force and then transmitted to the robotic arm controller for execution. The
218 grasping action is carried out by a force-controlled gripper, which maintains
219 a stable hold while preventing excessive compression or surface damage, ult-
220 imately achieving reliable and damage-free grasping.

221 By combining the commonsense reasoning capabilities of language models
222 with the geometric perception capabilities of point cloud and image data, this
223 framework offers an adaptable and generalizable solution for damage-free grasp-
224 ing across a wide variety of produce types in flexible agricultural robotics.

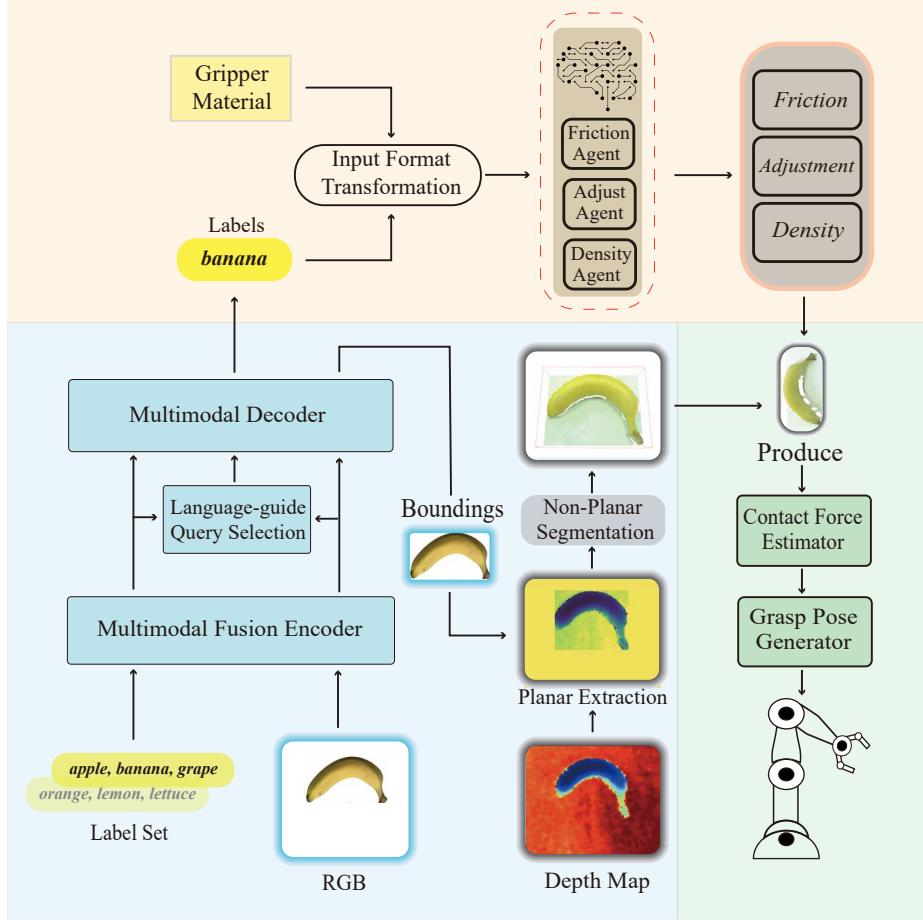


Fig. 1. Architecture overview of the proposed framework, composed of three main modules, (1) 3D perception module in blue, (2) LLM-based commonsense physical reasoning in orange, (3) grasp execution in green.

225 3.2. 3D Perception Module

226 Accurate 3D perception is essential for reliable produce localization and ma-
227 nipulation in unstructured agricultural environments. To this end, the per-
228 ception module integrates both semantic and geometric cues to achieve robust
229 detection and volume estimation. As illustrated in the following subsections,
230 we first employ a visual-language-based detection pipeline to localize target
231 produce using RGB images and category prompts. Subsequently, we introduce
232 a depth-enhanced, geometry-aware estimation process that leverages masked
233 point clouds and commonsense reasoning to infer volumetric properties of the
234 localized produce.

235 3.2.1. *Multimodal Instance Detection*

236 Grounding DINO is an open-vocabulary object detection model derived from
237 the DETR family, enabling flexible detection guided by textual prompts. Its
238 architecture comprises four primary components: an Image Backbone, a Text
239 Backbone, a Feature Enhancer that employs a dual-stream Transformer to fuse
240 and enhance visual and textual features, and the core modules, LGQS and
241 Cross-Modality Decoder.

242 As shown in Fig. 2, to facilitate effective multimodal interaction, our frame-
243 work introduces a multimodal fusion encoder. The core function of the mul-
244 timodal fusion encoder is to fuse image and text features. Initially, text and
245 image features are processed through a feature extractor, after which they enter
246 the cross-modality attention layer to enhance the interaction between the two
247 modalities. Following that, self-attention and deformable self-attention layers
248 are applied to further optimize the feature alignment, ensuring that the model
249 can extract and integrate effective information from different modalities, provid-
250 ing strong feature representations for downstream tasks. This process ensures
251 efficient multimodal information fusion, allowing the model to make accurate
252 judgments in complex scenarios.

253 The multimodal query-based decoder then takes the aligned features and
254 generates 2D coordinates of the target produce, enabling precise localization.
255 The decoder further optimizes the aligned features and assists in semantic label
256 alignment, ensuring the accuracy and reliability of the feature representation.
257 In addition, the decoder is not just a feature processor; it also integrates seman-
258 tic and geometric information, providing more accurate input for downstream
259 tasks, such as grasp pose estimation. With this powerful encoder-decoder ar-
260 chitecture, our framework significantly enhances grasp prediction for irregular
261 produce categories, ensuring both accuracy and stability in grasp pose estima-
262 tion.

263 Following multimodal feature fusion, the LGQS module plays a pivotal role
264 in generating semantically aligned queries for object detection. Such module
265 generates semantically guided Cross-Modality Queries by combining image fea-
266 tures with language prompts. These queries are then passed to the decoder
267 to predict the object Bounding Box. This design enables the model to align
268 semantics and localize regions based on arbitrary natural language descriptions,

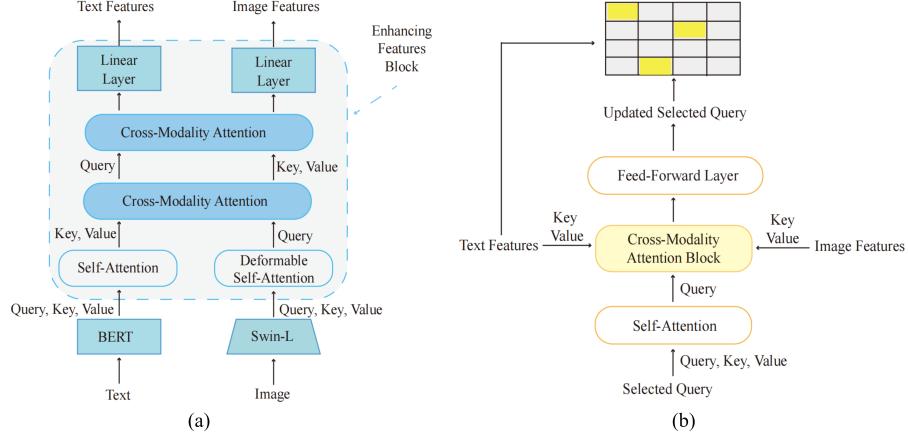


Fig. 2. Illustration of the encoder and decoder. (a) shows the Multimodal Feature Fusion Encoder, and (b) shows the multimodal query-based decoder.

equipping it with zero-shot detection capabilities (Hu et al., 2025b).

While the encoder-decoder modules are responsible for feature alignment and spatial localization, the LGQS module plays a critical role in bridging semantics with vision. Its internal computation process is illustrated in Fig. 3. The module takes as input the image features and text features produced by the respective backbones and computes a cross-modal similarity matrix, where each element reflects the semantic affinity between an image token and a text token. To identify the most semantically relevant regions in the image, the highest relevance score of each image token with respect to the entire text prompt is calculated, producing a relevance vector. Based on this vector, the image tokens with the highest semantic scores, corresponding to the most text-aligned visual regions, are selected and used as input queries for the cross-modality decoder.

In our system, we inherit the core architecture of Grounding DINO and adopt its multimodal fusion encoder and LGQS module to produce detection based on RGB images and a predefined label set. The LGQS module computes a similarity matrix between image and text features and selects the most relevant image tokens by identifying the maximum semantic response. The index selection process can be formally expressed as follows:

$$\mathbf{I}_{N_q} = \text{Top}_{N_q} \left(\text{Max}^{(-1)} (\mathbf{X}_I \mathbf{X}_T^\top) \right). \quad (1)$$

where Top picks the most relevant image tokens, N_q denotes the number of LGQS module output indices fixed at 900, $\text{Max}^{(-1)}$ extracts each image token's maximum relevance to the text, $\mathbf{X}_I \in \mathbb{R}^{N_I \times d}$ is the image features, $\mathbf{X}_T \in \mathbb{R}^{N_T \times d}$ is the text features, N_I and N_T denote the number of image and text tokens, and d is the shared embedding dimension.

To improve the system's flexibility and adaptability, a user-defined produce list mechanism is introduced, enabling dynamic adjustment of detection cat-

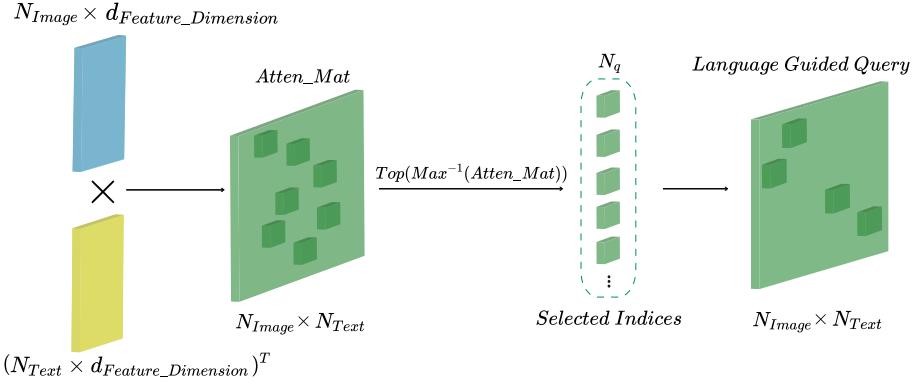


Fig. 3. Brief schematic of the LGQS module, where image and text features are combined to form an attention matrix, and the image token indices with the highest semantical score are selected to generate language-guided queries.

294 egories in response to specific grasping requirements. For example, produce
 295 categories such as blueberries or spinach can be dynamically added, and irrele-
 296 vant categories can be filtered out based on specific task requirements. This
 297 design improves the system’s practical applicability and robustness in dynamic
 298 agricultural environments.

299 *3.2.2. Geometry-aware Volume Estimation*

300 To improve volume estimation accuracy and enhance system adaptability,
 301 this study introduces a multi-stage estimation pipeline that combines visual
 302 perception with LLM-based commonsense reasoning. As shown in Fig. 4, the
 303 pipeline first acquires RGB images and the corresponding 2D coordinators of
 304 the target produce. Based on these coordinators, a 2D masked image is gen-
 305 erated and aligned with the depth map to isolate the region of interest, which
 306 is then transformed into point cloud data using the camera intrinsic parameter
 307 matrix. Planar Segmentation removes background surfaces, and the remaining
 308 points undergo Point Cloud Partitioning and Outlier Removal before Bound-
 309 ing Box Extraction provides a coarse geometric volume estimate. The original
 310 depth map is progressively refined through masking, point cloud filtering, and
 311 geometric fitting, effectively isolating the target produce, removing background
 312 interference, and constructing an enclosing bounding box. While this geometric
 313 enclosure often leads to overestimation due to redundant space, these visual
 314 results clearly illustrate how geometric information is extracted and utilized in
 315 the volume estimation process.

316 Regarding prompt design and further data processing, we observed that
 317 using more concrete and descriptive natural language expressions, for example
 318 "a piece of lettuce" outperformed abstract labels like "lettuce" by providing
 319 clearer semantic cues. This improvement enabled precise localization of target
 320 produce within cluttered scenes and laid a more robust visual groundwork for

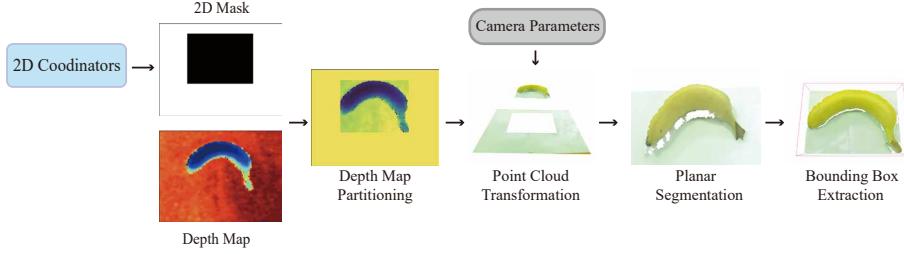


Fig. 4. Multi-stage pipeline for commonsense reasoning, combining 2D masks generation and depth maps converted to point clouds, followed by filtering and bounding box extraction.

321 downstream grasp parameter reasoning. In light of Grounding DINO’s image-
 322 text pair input format, we predefine a set of common produce categories as a
 323 default token pool for category matching during inference. With this design,
 324 the module requires only a single RGB image to automatically perform produce
 325 detection, without the need for additional textual input.

326 Once the target region is localized, its 2D mask is aligned with the corre-
 327 sponding depth map to isolate the produce’s depth profile. This masked depth
 328 data is then converted into a partial point cloud using the camera’s intrin-
 329 sic parameters. Subsequent processing includes planar segmentation to remove
 330 background surfaces, Euclidean clustering to separate the target from nearby
 331 objects, and statistical outlier removal to reduce noise. The cleaned point cloud
 332 is then fitted with an oriented bounding box to obtain coarse geometric mea-
 333 surements, which serve as the basis for initial volume estimation.

334 Together, these two modules provide a complete 3D perception pipeline that
 335 enables both robust instance detection and geometry-aware volume estimation,
 336 forming the basis for downstream grasping and manipulation.

337 3.3. LLM-Based Reasoning Module

338 The module first receives the semantic label of the target produce from
 339 the 3D perception module, such as banana, together with the gripper material
 340 specification. These inputs are processed by an Input Format Transformation
 341 stage, which re-formats the information for reasoning by a set of LLMs-driven
 342 agents, including FrictionAgent, AdjustAgent, and DensityAgent.

343 Without directly accessing point cloud data, the AdjustAgent leverages the
 344 semantic label to infer a canonical geometric primitive, such as an ellipsoid,
 345 cone, or cylinder-based on shape priors learned during pretraining. From this
 346 inferred geometry, the agent estimates the occupancy ratio within the bounding
 347 box, yielding a volume correction factor, which is used to refine the coarse
 348 volume estimate resulting in a more accurate produce volume. DensityAgent
 349 and FrictionAgent infer the produce density based on semantic knowledge and
 350 prior learning and estimate the friction coefficient between the produce surface
 351 and the gripper contact material, respectively.

Finally, the module further leverages LLMs to infer the produce's density and friction coefficient. These parameters, combined with Coulomb's law of friction, are used to calculate the minimum grasping force required for stable manipulation. First, the formula used to calculate the mass of tacet produce is as follows:

$$m = \rho \cdot k \cdot V \quad (2)$$

where ρ is the produce's density, k is the volume correction factor, and V is the estimated coarse volume. Then, the grasp force calculation formula is as follows:

$$F_{\text{grasp}} = \mu \cdot m \cdot g \quad (3)$$

where F_{grasp} is the required grasping force, μ is the friction coefficient and g is the gravitational acceleration valued as 9.81 m/s^2 . The estimated force is finally passed as a control parameter to the force-controlled gripper module.

By incorporating high-level semantic reasoning and physical property inference, the proposed framework substantially improves the accuracy of volume estimation and grasping force prediction. This enables the force-controlled gripper to execute stable, damage-free manipulation of complex, irregular, or deformable produce.

3.4. Grasping Execution Module

VMGNet is a grasp representation framework that generates 2D planar grasping predictions from RGB or RGB-D images. Its core architecture comprises the Visual State Space (VSS) and the Fusion Bridge Module (FBM), which together enable efficient grasp pose inference and planning. The VSS component compresses image information to reduce computational complexity, while the FBM enhances the model's ability to capture multi-scale features, thereby improving the accuracy and robustness of grasp point generation (Jin et al., 2024).

Although VMGNet exhibits strong grasp representation capabilities under 2D image inputs, its original design primarily targets structurally stable general objects. This makes it less adaptable to agricultural scenarios, where produce often exhibits soft textures, irregular morphologies, and unstructured features. To enhance VMGNet's performance in damage-free agricultural grasping tasks, we implement targeted optimizations to its input pipeline and grasping strategy.

In the case of VMGNet, the Region of Interest (ROI) is directly defined as a graspable region and parameterized using a grasp rectangle representation, expressed as follows:

$$ROI = (x, y, \theta, h, w) \quad (4)$$

where (x, y) denotes the center coordinates of the grasp rectangle, θ represents the grasp angle relative to the horizontal axis, and h, w correspond to the height and width of the proposed grasp. These parameters are regressed directly from fused multi-scale features via the decoder, without requiring explicit cross-modal alignment or segmentation.

391 We introduce the LGQS module from Grounding DINO to extract the se-
 392 mantic label of the target produce and generate a precise 2D mask for localizing
 393 its boundary within the image. This semantic-guided mechanism ensures that
 394 the detection region is both accurate and complete, mitigating grasp errors
 395 caused by false detections or background interference. We then align the de-
 396 tected produce region with its corresponding depth map to construct a fused
 397 image input focused on the target produce, significantly improving the quality
 398 of grasp candidate regions and reducing interference from non-target areas. To
 399 demonstrate the practical effectiveness of VMGNet in real-world agricultural
 400 scenarios, we visualize the predicted grasp poses for various delicate produce.
 401 As shown in Fig. 5, the VMGNet framework accurately identifies ROI for grasp-
 402 ing, as depicted by the bounding boxes on the produce. These bounding boxes
 403 guide the Contact Force Estimator to predict stable grasp forces by ensuring
 404 that the contact force is appropriate for each object’s physical characteristics,
 405 thereby preventing damage during manipulation.

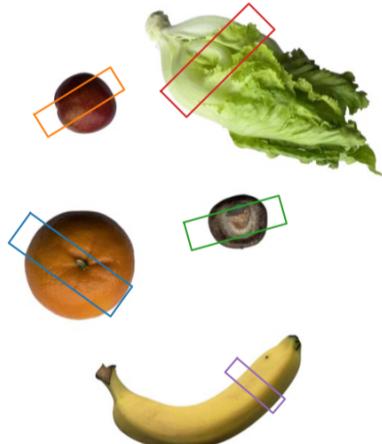


Fig. 5. Graspable regions predicted by VMGNet for various agricultural produce for minimal manipulation damage.

406 Considering the sensitivity of force control in handling delicate produce, we
 407 design a Contact Force Estimator based on LLM-inferred physical attributes.
 408 This module utilizes key parameters, such as density, friction coefficient, and a
 409 volume correction factor inferred through language model reasoning, to estimate
 410 the minimum stable contact force required for damage-free grasping. The esti-
 411 mated force guides the grasp planning module in determining whether a grasp
 412 action complies with damage-free constraints, thereby preventing produce dam-
 413 age caused by excessive gripping force or inaccurate object understanding.

414 **4. Experiment**

415 This experiment aims to evaluate the proposed multimodal damage-free
416 grasping framework in terms of LLMs reasoning-based performance and grasp
417 force control, with a focus on comparing different estimation pathways and
418 assessing the effectiveness of LLMs in physical property reasoning and grasp
419 feasibility.

420 **4.1. Experimental Setup and Data Collection**

421 We implement the proposed framework following the architecture described
422 in [Section 3](#), which integrates a 3D perception module, an LLM-based physical
423 reasoning module, and a grasp execution module. We conduct all experiments
424 using four NVIDIA RTX 4090 GPUs with each 24GB of memory. The imple-
425 mentation is based on PyTorch 2.0.1 and Python 3.8.20. All text embedding
426 processing and tokenization steps are performed using BERT. The training and
427 evaluation are carried out on Ubuntu 20.04.

428 To enhance the detection of agricultural produce categories insufficiently
429 represented in the pretrained model, we fine-tuned Grounding DINO on an
430 augmented dataset combining Objects365 V1 ([Shao et al., 2019](#)), OpenImages
431 V6 ([Kuznetsova et al., 2020](#)), and GoldG([Kebe et al., 2021](#)). RGB and depth
432 data are captured using an Intel RealSense D435 camera, producing RGB-D
433 images and point clouds for 3D bounding box construction and mass estimation.
434 For each produce sample, high-precision electronic scales are used to obtain
435 ground-truth mass measurements.

436 The experimental dataset consists of 25 typical produce categories, covering
437 berries, fruits, mushrooms, leafy vegetables, hollow produce, and part of the
438 produce structure. These samples exhibit diverse shapes, densities, surface tex-
439 tures, and structural stiffness, reflecting the challenges of real-world agricultural
440 grasping. Representative examples are shown in [Fig. 6](#), and the full category
441 list is provided in the Appendix.

442 **4.2. Evaluation Metrics**

443 To comprehensively evaluate the effectiveness of the proposed multimodal
444 damage-free grasping framework, we design three experimental tasks: compara-
445 tive reasoning performance analysis across different LLMs, mass estimation ac-
446 curacy assessment, and grasp force feasibility validation, each with well-defined
447 quantitative metrics.

448 **4.2.1. Reasoning-based Tasks Metrics**

449 To ensure consistency and comparability across reasoning-based tasks, we
450 adopt a unified set of evaluation metrics for both the mass estimation task and
451 the LLMs reasoning capability comparison task. Specifically, we assess model
452 performance along two dimensions: accuracy and stability.

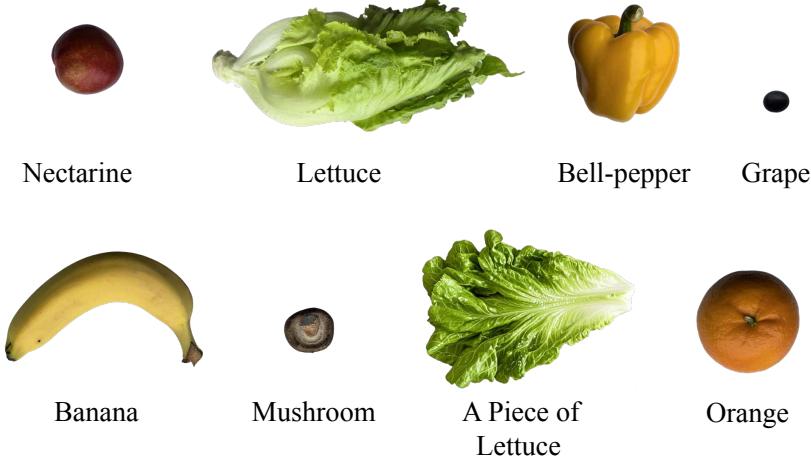


Fig. 6. Representative samples from the experimental dataset, covering diverse produce categories with varying shapes, textures, and physical properties.

453 Accuracy is quantified using the Mean Absolute Percentage Error (MAPE)
 454 as follows:

$$\text{MAPE} = \frac{1}{N} \sum_{i=1}^N \left(\left| \frac{\hat{y}_i - y_i}{y_i} \right| \right) \times 100\% \quad (5)$$

455 where \hat{y}_i is the predicted value of the i -th sample, y_i is the ground truth value of
 456 the i -th sample, and N represents the total number of samples. Lower MAPE
 457 indicates higher estimation accuracy. Stability of predictions across samples is
 458 measured by the Standard Deviation (STD) as follows:

$$STD = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2} \quad (6)$$

459 where x_i is the i -th observed value, \bar{x} is the sample mean, N is the total number
 460 of samples. Smaller STD reflects more consistent predictions.

461 The ground-truth mass of each sample is measured using a high-precision
 462 electronic scale, serving as the baseline for evaluation. Estimation accuracy is
 463 quantified by the MAPE, where lower values indicate higher accuracy. To ensure
 464 numerical stability, a minimum denominator threshold is applied in the MAPE
 465 computation to avoid division by very small ground-truth values. Prediction
 466 stability is assessed by the STD of the errors.

467 4.2.2. Grasp Feasibility Evaluation Metrics

468 To assess the accuracy of predicted grasping forces in a safety-aware manner,
 469 we adopt a tolerance-based hit rate metric. For each produce category c , we
 470 denote the reference *success force* of sample i as $F^{(i)}$. Based on empirical

471 fragility considerations, an asymmetric tolerance interval is defined as

$$[L^{(i)}, U^{(i)}] = [(1 - \delta_c^-)F^{(i)}, (1 + \delta_c^+)F^{(i)}], \quad (7)$$

472 where δ_c^- and δ_c^+ represent the category-specific lower and upper tolerance ratios,
473 respectively. Given a calibrated model prediction $\hat{F}^{(i)}$, a hit indicator is
474 defined as

$$h^{(i)} = \begin{cases} 1, & \text{if } L^{(i)} \leq \hat{F}^{(i)} \leq U^{(i)}, \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

475 The **hit rate** for model m within category c is then computed as the empirical
476 average:

$$\text{HitRate}_{m,c} = \frac{1}{N_c} \sum_{i=1}^{N_c} h^{(i)}, \quad (9)$$

477 where N_c denotes the number of samples in category c . This metric reflects the
478 proportion of predictions that fall into the physically safe and effective grasping
479 force interval. It complements continuous error metrics (e.g., MAPE) by directly
480 quantifying the likelihood that a model's predictions are simultaneously non-
481 damaging and sufficient for successful grasping.

482 Together, these metrics provide a comprehensive and task-relevant evaluation
483 of both the physical property reasoning and the downstream force-controlled
484 grasp execution capabilities of the proposed system, reflecting its practicality
485 and reliability in real-world agricultural environments.

486 4.3. Perception-Fusion Mass Estimation

487 To systematically evaluate the contribution of geometric perception to mass
488 estimation, we compare two paradigms: (i) a **reasoning-only** setting, where the
489 LLMs directly infer the mass of the target produce in kilograms from its semantic
490 label without any geometric perception, and (ii) our proposed **perception-fusion**
491 framework, which explicitly incorporates 3D information. In the latter,
492 the module first derives a 3D bounding box and corresponding volumetric estimate
493 from RGB-D observations, then applies an LLM-inferred shape correction
494 factor to refine the raw volume, and subsequently combines the corrected volume
495 with an LLM-inferred density to obtain the final mass prediction. This
496 design enables a clear isolation of the benefits brought by integrating geometric
497 structure into the reasoning process.

498 Among the models assessed, such as Deepseek-r1, Deepseek-v3, qwen-plus,
499 QWen-2.5 (14B), QWen-max, ChatGPT-4o, and Claude-3.5-Sonnet, the 2.5
500 (14B) model represents the smallest publicly released system available at the
501 time of study. While its limited parameter count constrains baseline reasoning
502 capability, we mitigate this by performing domain-specific fine-tuning on
503 the Camel dataset (Li et al., 2023a). This specialized adaptation allows the
504 model to acquire inductive biases tailored for volumetric and density reasoning,
505 forming the backbone of our final framework. The results highlight not
506 only the value of geometric integration but also the potential of compact LLMs

507 when reinforced through carefully curated training, with experimental results
 508 across different LLMs summarized in [Table 1](#) and further details provided in
 509 the Appendix.

Table 1

Average MAPE and STD in produce mass estimation using reasoning-only vs.
 perception-fusion methods, with model parameter sizes (B), on a custom dataset of
 seven produce categories.

Models	Reasoning-only		Perception-fusion		Parameter (B)
	MAPE	STD	MAPE	STD	
ChatGPT4o	792.842	1022.937	88.349	95.469	220
Claude-3.5-Sonnet	142.843	2126.330	81.729	99.335	≥ 175
Deepseek R1	312.601	284.372	190.326	220.491	671
Deepseek V3	170.083	125.146	82.123	69.441	671
QWen-max	301.782	249.868	81.2889	61.204	≥ 72
QWen-plus	227.963	234.358	83.822	82.937	≥ 72
Ours	376.804	340.810	45.809	23.590	14

510 In the task of produce mass estimation, our experimental results system-
 511 atically demonstrate the pronounced advantages of incorporating point cloud
 512 assistance. First, in terms of overall error levels, the introduction of geometric
 513 information from point clouds significantly improves both prediction accuracy
 514 (MAPE) and stability (STD) across almost all evaluated models. For instance,
 515 ChatGPT4o under the reasoning-only setting exhibits an average MAPE as high
 516 as 792.84% with an STD exceeding 1000, indicating severe instability when rely-
 517 ing solely on semantic inference across diverse produce categories. In contrast,
 518 under the perception-fusion setting, its MAPE drops to 88.35% and STD to
 519 95.47, a reduction of nearly 90%, clearly demonstrating that three-dimensional
 520 geometric cues effectively compensate for the inherent limitations of purely se-
 521 mantic reasoning. Similarly, the MAPE of Claude-3.5-Sonnet decreases from
 522 142.84% to 81.73%, highlighting that even large-scale models benefit substan-
 523 tially from the integration of geometric perception.

524 Second, the advantages of point cloud integration are even more pronounced
 525 for medium- and small-scale models. For example, Qwen-plus and Qwen-max
 526 under the reasoning-only paradigm exhibit average errors in the 200%-300%
 527 range, but these are reduced to 83.82% and 81.29%, respectively, under the
 528 perception-fusion setting, corresponding to relative improvements of more than
 529 60%. This demonstrates that geometric structural information serves as a highly
 530 effective complement for models with limited parameter scales, enabling them
 531 to achieve performance levels comparable to or even surpassing those of much
 532 larger models. Of particular note is our 14B-parameter model (Ours), which
 533 achieves the lowest overall error with the lowest MAPE(45.81%) and the small-
 534 est STD (23.59) under the perception-fusion setting. This performance not
 535 only far exceeds that of its reasoning-only counterpart, which leads to relatively

536 high MAPE(3776.804%) and STD(340.810), but also substantially outperforms
537 larger models such as ChatGPT4o and Claude-3.5-Sonnet. These findings un-
538 derscore that through the synergy of geometric perception and semantic rea-
539 soning, even medium-scale models can deliver performance superior to state-of-
540 the-art large-scale models in complex produce mass estimation tasks.

541 More importantly, from an agricultural application perspective, different
542 produce categories exhibit substantial variations in geometric morphology and
543 physical properties, and the point cloud-assisted method demonstrates strong
544 stability across these heterogeneous classes. For instance, in categories such as
545 leafy vegetables and parts of produce structures, reasoning-only methods often
546 suffer from large fluctuations in estimation errors, whereas the perception-fusion
547 strategy effectively reduces both bias and variance, thereby substantially low-
548 ering the risk of damage in non-destructive grasping. For fragile produce such
549 as berries and mushrooms, point cloud perception not only improves the pre-
550 cision of volumetric estimation but also provides more reliable priors to guide
551 subsequent grasping strategies. This cross-category robustness carries direct
552 significance for automated harvesting, sorting, and packaging in agricultural
553 production. Detailed per-category statistics supporting these observations are
554 provided in the Appendix.

555 4.4. Reasoning Performance Comparison Across LLMs

556 In damage-free grasping tasks, accurate force control strategies critically de-
557 pend on reliable estimation of the target object’s mass and frictional properties.
558 As mass is not directly measurable during inference, commonsense reasoning in-
559 formed by visual cues and semantic labels becomes essential for approximating
560 physical properties. Therefore, this section focuses on comparing the perfor-
561 mance of different LLMs in commonsense reasoning and evaluating their ability
562 to support grasp control strategies.

563 This study selects seven state-of-the-art LLMs, including GPT-4o, Qwen-
564 plus, Qwen-max, Qwen-2.5 with fine-tuned (Ours), Deepseek-r1, Deepseek-v3,
565 and Claude-3.5-Sonnet. Each model receives standardized semantic labels of
566 produce as input and outputs both density and volume correction factors. These
567 values are combined with volume estimates obtained from RGB-D point cloud-
568 derived 3D bounding boxes, enabling mass calculation. Based on the inferred
569 mass and friction coefficients, the system applies Coulomb’s friction model to
570 calculate the theoretical minimum grasping force. This force is then applied
571 using a force-controlled gripper equipped with closed-loop control capabilities.
572 The system records whether the grasp is successful and whether any deformation
573 or damage occurred.

574 In parallel, the MAPE is used to quantitatively compare the mass estima-
575 tion accuracy of the seven LLMs across seven representative produce structural
576 categories. The results are summarized in [Table 2](#).

577 Among all evaluated models under the perception-fusion setting, Ours achieves
578 the lowest overall average MAPE of 45.81%, substantially outperforming larger-
579 scale models such as ChatGPT-4o (88.35%) and Claude-3.5-Sonnet (81.73%).

Table 2

Comparison of MAPE across LLMs on seven representative produce categories under point cloud-assisted mass estimation.

	MAPE						
	Deepseek		Qwen		ChatGPT	Claude	Ours
	r1	v3	max	plus	4o	3.5-Sonnet	
1	46.889	36.958	41.558	69.246	50.002	45.408	69.393
2	47.579	42.513	55.664	63.831	47.563	49.745	52.192
3	825.177	21.704	49.333	31.094	34.266	37.068	22.008
4	35.135	36.565	42.739	42.289	36.876	35.1	47.676
5	335.094	136.478	180.718	135.478	56.999	171.998	78.484
6	12.236	250.664	144.071	233.544	383.41	201.499	41.597
7	30.171	49.982	54.939	11.271	9.324	31.286	9.311
All	190.326	82.123	81.289	83.822	88.349	81.729	45.809

Note: 1 stands for Berries; 2 for Round Fruits; 3 for Elongated Fruits; 4 for Elliptical Fruits; 5 for Leafy Vegetables; 6 for Part of Produce Structure; and 7 for Mushroom Family.

580 This demonstrates its strong generalization capability and robustness, highlighting
 581 its suitability as a compact yet reliable inference backbone for multi-category
 582 produce grasping systems. In contrast, Deepseek-r1 exhibits the weakest gen-
 583 eralization with an average MAPE of 190.33%, indicating limited adaptability
 584 across structurally diverse produce.

585 At the category level, Ours consistently delivers leading performance, achiev-
 586 ing the lowest MAPE in Mushroom Family (9.31%) and Part of Produce Struc-
 587 ture (41.60%), both of which are characterized by fragility and structural com-
 588 plexity. Deepseek-v3 shows relative strength in Elongated Fruits (21.70%), sug-
 589 gesting an ability to capture density-to-volume relationships in slender geome-
 590 tries. However, the Deepseek series also reveals clear bottlenecks, for example,
 591 Deepseek-r1 reached 47.58% in Round Fruits and 335.09% in Leafy Vegeta-
 592 bles, underscoring its instability in categories that are either density-variable or
 593 structurally loose.

594 Other models demonstrate intermediate patterns. ChatGPT-4o exhibits bal-
 595 anced performance across most categories, reflecting a trade-off between seman-
 596 tic understanding and reasoning ability. Claude-3.5-Sonnet achieves strong re-
 597 sults in Mushroom Family (31.29%), but its performance fluctuated significantly
 598 in other categories, such as Leafy Vegetables (171.99%), indicating dependence
 599 on semantic clarity and training data relevance.

600 Overall, these results underscore the importance of integrating geometric
 601 perception with language-based reasoning in agricultural mass estimation. The
 602 cross-category variations in MAPE highlight that not all state-of-the-art LLMs
 603 generalize equally well. Notably, the fine-tuned compact model (Ours) demon-
 604 strates the greatest potential as a reliable backbone for damage-free robotic
 605 grasping, particularly in scenarios involving diverse, fragile, and structurally
 606 ambiguous produce.

607 4.5. Comparison with Baseline Approaches

608 To further evaluate the practical adaptability and control feasibility in real-
 609 world grasping scenarios, a grasp force validation experiment is conducted be-
 610 tween Ours and DeliGrasp (Xie et al., 2024). Evaluation metrics include the hit
 611 rate, which measures the proportion of predictions falling within asymmetric
 612 safety tolerances, and the MAPE, which quantifies continuous estimation error.
 613 Together, these results provide insight into how different baseline methods bal-
 614 ance accuracy and safety compliance in damage-free grasping, where a higher
 615 hit rate and lower MAPE indicate a better trade-off between the two.

616 A calibration factor serves as a lightweight linear adjustment that rescales
 617 the raw force estimates produced by each model. Different categories of produce
 618 often exhibit systematic overestimation or underestimation due to variations in
 619 density, geometry, or fragility, and a single global model may not capture these
 620 nuances. By learning a multiplicative factor s for each model-category pair, we
 621 align predictions with the physical reference scale, thereby reducing systematic
 622 bias while preserving relative variations across samples. Experimental results
 623 are illustrated in [Table 3](#).

Table 3

Performance comparison of DeliGrasp (Xie et al., 2024) and Ours across produce categories, reporting hit rate, MAPE, and calibration factor.

Categories	DeliGrasp (Xie et al., 2024)			Ours		
	Hit Rate	MAPE	Factor	Hit Rate	MAPE	Factor
1	0.667	0.147	0.756	0.8	0.136	0.78
2	0.6	0.354	0.661	0.8	0.352	0.7
3	0.667	0.311	0.673	1	0.206	0.585
4	0.375	0.251	0.554	0.667	0.378	0.98
5	0.375	0.412	0.56	0.25	0.76	0.85
6	0.889	0.165	0.708	0.667	0.242	0.6
7	0.667	0.424	0.743	0.667	0.659	0.89
All	0.606	0.295	-	0.693	0.39	-

Note: 1 stands for Berries; 2 for Round Fruits; 3 for Elongated Fruits; 4 for Elliptical Fruits; 5 for Leafy Vegetables; 6 for Part of Produce Structure; and 7 for Mushroom Family.

624 At the aggregate level, Ours achieves a superior balance between accuracy
 625 and safety compliance. The results show that its average hit rate reaches 0.693,
 626 whereas DeliGrasp attains only 0.606, demonstrating a more consistent ability
 627 to generate grasping force estimates that remain within the asymmetric safety
 628 tolerances. Although the MAPE is slightly higher at 0.390 compared with 0.295
 629 for DeliGrasp, the improvement in hit rate highlights the stronger robustness of
 630 Ours when applied to fragile produce. This indicates that in real-world scenarios,
 631 our framework is more effective at preventing damage caused by excessive
 632 force while still maintaining acceptable predictive precision.

633 A category-wise analysis further emphasizes these advantages. In highly
634 fragile categories such as Berries and Round Fruits, Ours attains hit rates of
635 0.8 in both cases, exceeding the 0.667 and 0.6 achieved by DeliGrasp, while also
636 delivering lower estimation errors. In Elongated Fruits, Ours achieves a perfect
637 hit rate of 1.0, in contrast to only two-thirds for DeliGrasp, and simultaneously
638 yields smaller errors. These results demonstrate that our framework generalizes
639 more effectively across produce with diverse geometries and structural proper-
640 ties, providing accurate estimates while preserving the structural integrity of
641 delicate produce.

642 The analysis of calibration factors provides additional evidence of robustness.
643 In several categories, the scaling factors learned by Ours remain close to unity,
644 with values such as 0.78 for Berries and 0.98 for Part of Produce Structure,
645 indicating minimal systematic bias between predicted and reference values. In
646 contrast, DeliGrasp exhibits much stronger deviations, with factors of only 0.554
647 for Part of Produce Structure and 0.56 for Leafy Vegetables, suggesting larger
648 inherent biases. The reduced reliance of Ours on calibration adjustments con-
649 firms the intrinsic stability of its raw predictions and enhances interpretability,
650 as the outputs remain well aligned with the physical reference scale.

651 Beyond the numerical results, a qualitative comparison of grasping outcomes
652 further underscores the effectiveness of our framework in achieving damage-free
653 manipulation. [Fig. 7](#) provides a side-by-side illustration of grasping instances
654 produced by our method and DeliGrasp. In the top row, grasps executed by our
655 framework demonstrate stable force regulation that firmly secures the objects
656 while preserving their original structural integrity. The absence of deformation
657 or surface damage highlights the capability of our method to deliver reliable
658 control even for fragile and irregularly shaped produce. In contrast, the bot-
659 tom row depicts grasps from DeliGrasp, where uneven or excessive forces often
660 result in visible deformation, slippage, or local crushing. These shortcomings
661 are particularly evident in delicate categories such as leafy vegetables and elon-
662 gated fruits, where excessive compression compromises the objects' usability
663 and market value.

664 The visual evidence aligns closely with the quantitative analysis reported
665 in [Table 3](#). Categories such as berries, round fruits, and elongated fruits show
666 both higher hit rates and lower estimation errors for our method, which directly
667 translates into the more consistent and damage-free grasps visible in [Fig. 7](#).
668 DeliGrasp, despite achieving a lower average MAPE, fails to guarantee force
669 regulation within the asymmetric safety tolerances, which manifests as the phys-
670 ical deformations shown in the figure. This illustrates that a lower error metric
671 alone is insufficient for safe grasping, and that the ability to maintain forces
672 within the safe operational range is more critical for real-world applications.

673 Taken together, these results demonstrate that our method offers not only
674 quantitative improvements but also practical advantages in terms of adaptabil-
675 ity, robustness, and safety compliance. By achieving accurate force predictions
676 that directly translate into stable and damage-free manipulation, our framework
677 provides a stronger foundation for damage-free grasping in real-world agricul-
678 tural and industrial environments.

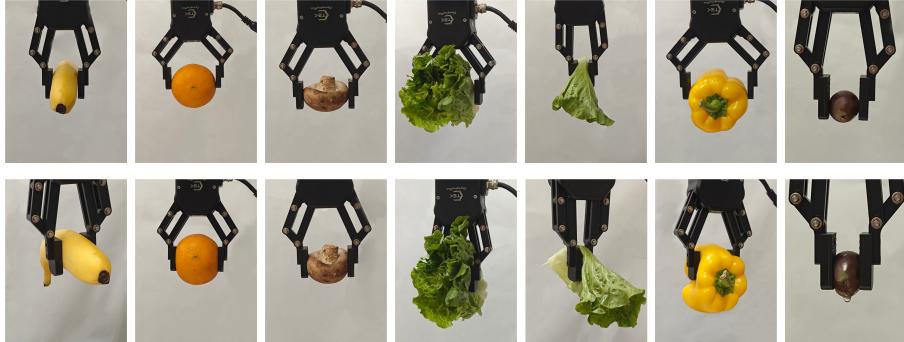


Fig. 7. Comparison of grasping performance. The top row shows successful grasping by Ours, while the bottom row displays results from DeliGrasp (Xie et al., 2024), which may lead to deformation or even damage.

679 5. Conclusion

680 We propose a novel grasping framework, which builds upon the planar grasping
 681 framework by incorporating semantic enhancement and commonsense physical
 682 reasoning to improve adaptability and robustness across diverse produce
 683 categories and morphologies. The framework integrates the LGQS module from
 684 Grounding DINO, enabling deep alignment between natural language inputs
 685 and visual features. This enables the framework not only to localize poten-
 686 tial grasp regions but also to identify the produce category and key attributes,
 687 thereby supporting semantically driven grasp strategy generation. Building on
 688 semantic identification, further leverages the extensive knowledge base and rea-
 689 soning capability of LLMs. By interpreting class-specific semantic descriptions
 690 of produce, the framework can autonomously infer implicit physical attributes
 691 such as density and friction coefficient, which are then used to estimate appro-
 692 priate grasping forces. This semantic-physical reasoning mechanism provides a
 693 robust cognitive foundation for executing damage-free grasping tasks. In addi-
 694 tion, this framework utilizes point cloud data of the target produce, combined
 695 with the heuristic estimation ability of LLMs, to improve volume estimation and
 696 grasp region identification for irregularly shaped produce such as density-varying
 697 mushrooms, curved bananas, or leafy greens with internal voids. This enhances
 698 generalization and adaptability in handling complex produce geometries.

699 Despite the promising results achieved in damage-free grasping of delicate
 700 produce, the framework still faces several limitations. First, LLMs exhibit insta-
 701 bility in density inference for structurally complex produce such as dragon-fruit
 702 and avocados. Second, the current volume correction mechanism relies on static
 703 factors rather than a learned model, limiting its accuracy. Future work will fo-
 704 cus on improving system performance by incorporating visual context cues to
 705 enhance the physical reasoning capabilities of LLMs and introducing a learn-
 706 able volume correction network to improve the geometric estimation accuracy
 707 for diverse produce types.

708 **CRediT authorship contribution statement**

709 **Ziye Zhang:** Conceptualization, Methodology, Software, Formal analysis,
710 Data Curation, Writing - Original Draft & Review & Editing, Visualization. **Xi-**
711 **aoyu Xia:** Conceptualization, Methodology, Software, Formal analysis, Data
712 Curation, Writing - Original Draft & Review & Editing, Visualization. **Yuhao**
713 **Jin:** Conceptualization, Methodology, Resources, Writing - Review & Edit-
714 ing. **Qizhong Gao:** Visualization, Writing - Review & Editing. **Lin Qiao:**
715 Resources, Formal analysis. **Jinglei Chen:** Formal analysis, Data Curation.
716 **Yong Yue:** Supervision, Funding acquisition, Validation. **Xiaohui Zhu:**

717 **Declaration of competing interest**

718 The authors declare that they have no known competing financial or personal
719 relationships that could have appeared to influence the work reported in this
720 paper.

721 **Acknowledgments**

722 The authors acknowledge the technical support and device contribution to
723 Xi'an Jiaotong-Liverpool University (XJTLU).

724 **Funding**

725 This research was funded by Suzhou Municipal Key Laboratory for Intelli-
726 gent Virtual Engineering (Szs2022004), the XJTLU Artificial Intelligence (AI)
727 University Research Centre, Jiangsu Province Engineering Research Centre of
728 Data Science and Cognitive Computation at XJTLU, and Suzhou Industrial
729 Park (SIP) AI innovation platform (YZCXPT2022103).

730 **Appendix A. Produce sample with corresponding category**

731 In [Table A1](#), produce items are classified by category and grasp-related char-
732 acteristics, providing a structured basis for analyzing grasp difficulty and esti-
733 mation challenges.

734 **Appendix B. Complete Experimental Data**

735 This appendix provides the complete per-category results of mass estimation
736 across seven representative produce categories using seven different LLMs. The
737 [Table B1](#) and [Table B2](#) report MAPE and STD under both reasoning-only and
738 perception-fusion settings, complementing the aggregated results presented in
739 the main text.

Table A1
Produce classification and their grasping-related characteristics.

Produce	Categories	Characteristic
grape, cherry, berry, strawberry, cherry-tomato	Berries	Small volume; Smooth surface; Risky of slippage
apple, tomato, nectarine	Round Fruits	Spherical; high density; Easy to predict
banana, cucumber	Elongated Fruits	Elongated and curved; Difficult to shape fit
lemon, avocado, orange, dragon-fruit, kiwi-fruit	Elliptical Fruits	Rough surface; Medium volume; Easy to deformation
bell-pepper	Hollow Produce	Hard to estimate
spinach, lettuce, Chinese-cabbage, baby-cabbage	Leafy Vegetables	Lightweight; Large volume; Loosely structured
watermelon-flesh, broccoli, a-piece-of-lettuce	Part of Produce Structure	Difficult to quantify size; Incomparable qualitatively
button-mushroom, shiitake-mushroom, King-Oyster-mushroom	Mushroom Family	Volume - Mass imbalance

Table B1

Complete per-category results of produce mass estimation across seven LLMs and seven representative categories. Reported are MAPE and STD under reasoning-only and perception-fusion paradigms.

Models	Categories	Reasoning-only		Perception-fusion	
		MAPE	STD	MAPE	STD
Deepseek R1	1	53.17	21.057	46.889	28.739
	2	35.543	5.009	47.579	12.827
	3	389.258	525.284	825.177	1140
	4	30.089	17.029	35.135	24.269
	5	462.955	703.512	335.094	309.018
	6	943.952	451.066	12.236	9.623
	7	273.237	267.648	30.171	18.964
	All	312.601	284.372	190.326	220.491
Deepseek V3	1	88.557	8.506	36.958	31.825
	2	359.015	578.433	42.513	29.127
	3	28.802	3.157	21.704	8.158
	4	24.862	21.976	36.565	15.733
	5	178.343	123.152	136.478	185.692
	6	448.224	106.216	250.664	209.68
	7	62.779	34.581	49.982	5.87
	All	170.083	125.146	82.123	69.441
QWen-max	1	85.084	12.012	41.558	34.592
	2	16.279	19.613	55.664	19.738
	3	7.972	5.74	49.333	0.146
	4	50.747	41.75	42.739	16.948
	5	1255.166	1053.88	180.718	145.58
	6	485.95	390.708	144.071	191.636
	7	211.276	225.374	54.939	19.789
	All	301.782	249.868	81.2889	61.204
QWen-plus	1	43.68	46.944	69.246	38.27
	2	35.717	44.879	63.831	9.007
	3	36.53	39.434	31.094	9.8
	4	38.085	19.375	42.289	19.308
	5	552.95	512.414	135.478	164.507
	6	830.379	916.582	233.544	322.436
	7	58.4	60.881	11.271	17.232
	All	227.963	234.358	83.822	82.937
ChatGPT4o	1	65.144	21.403	50.002	24.634
	2	36.828	10.438	47.563	22.433
	3	5.686	2.867	34.266	26.382
	4	51.503	32.251	36.876	13.642
	5	394.81	621.336	56.999	56.794
	6	4947.936	6463.047	383.41	521.184
	7	47.986	9.22	9.324	3.214
	All	792.842	1022.937	88.349	95.469

Note: 1 stands for Berries; 2 for Round Fruits; 3 for Elongated Fruits; 4 for Elliptical Fruits; 5 for Leafy Vegetables; 6 for Part of Produce Structure; and 7 for Mushroom Family.

Table B2

Complete per-category results of produce mass estimation across seven LLMs and seven representative categories. Reported are MAPE and STD under reasoning-only and perception-fusion paradigms (continue).

Models	Categories	Reasoning-only		Perception-fusion	
		MAPE	STD	MAPE	STD
Claude-3.5-Sonnet	1	137.061	196.967	45.408	23.365
	2	59.42	49.95	49.745	15.541
	3	159.26	38.671	37.068	4.514
	4	77.929	81.73	35.1	12.788
	5	237.538	10186.482	171.998	450.291
	6	265.742	4324.687	201.499	173.815
	7	62.948	5.825	31.286	15.029
	All	142.843	2126.330	81.729	99.335
Ours	1	79.356	28.154	69.393	34.377
	2	34.704	40.799	52.192	17.355
	3	12.624	5.555	22.008	2.33
	4	63.54	50.884	47.676	19.281
	5	1099.925	1021.552	78.484	51.089
	6	861.828	1171.669	41.597	37.315
	7	485.651	67.059	9.311	3.385
	All	376.804	340.810	45.809	23.590

Note: 1 stands for Berries; 2 for Round Fruits; 3 for Elongated Fruits; 4 for Elliptical Fruits; 5 for Leafy Vegetables; 6 for Part of Produce Structure; and 7 for Mushroom Family.

740 **References**

- 741 Ahn, M., Brohan, A., Brown, N., Chebotar, Y., Cortes, O., David, B., Finn,
742 C., Fu, C., Gopalakrishnan, K., Hausman, K., et al., 2022. Do as i can,
743 not as i say: Grounding language in robotic affordances. arXiv preprint
744 arXiv:2204.01691 doi:[10.48550/arXiv.2204.01691](https://doi.org/10.48550/arXiv.2204.01691).
- 745 Chen, K., Li, T., Yan, T., Xie, F., Feng, Q., Zhu, Q., Zhao, C., 2022. A soft
746 gripper design for apple harvesting with force feedback and fruit slip detection.
747 Agriculture 12, 1802. doi:[10.3390/agriculture12111802](https://doi.org/10.3390/agriculture12111802).
- 748 Chen, X., Sun, Y., Zhang, Q., Dai, X., Tian, S., Guo, Y., 2024. Three-
749 dimension object detection and forward-looking control strategy for non-
750 destructive grasp of thin-skinned fruits. Applied Soft Computing 150, 111082.
751 doi:[10.1016/j.asoc.2023.111082](https://doi.org/10.1016/j.asoc.2023.111082).
- 752 Chu, X., Deng, J., You, G., Liu, W., Li, X., Ji, J., Zhang, Y., 2025. GraspCOT:
753 Integrating physical property reasoning for 6-dof grasping under flexible lan-
754 guage instructions. arXiv preprint arXiv:2503.16013 doi:[10.48550/arXiv.2503.16013](https://doi.org/10.48550/arXiv.2503.16013).
- 755 Ghasemzadeh, P., Hempel, M., Wang, H., Sharif, H., 2023. Ggcnn: An
756 efficiency-maximizing gated graph convolutional neural network architecture
757 for automatic modulation identification. IEEE Transactions on Wireless Com-
758 munications 22, 6033–6047. doi:[10.1109/TWC.2023.3239311](https://doi.org/10.1109/TWC.2023.3239311).
- 759 Hu, F., Wu, F., Gu, H., Abbas, G., Alanazi, M.D., Othmen, S., Wang, J.,
760 Zhang, T., 2025a. Transforming agriculture with advanced robotic decision
761 systems via deep recurrent learning. Expert Systems with Applications 259,
762 125123. doi:[10.1016/j.eswa.2024.125123](https://doi.org/10.1016/j.eswa.2024.125123).
- 763 Hu, Z., Gao, K., Zhang, X., Yang, Z., Cai, M., Zhu, Z., Li, W., 2025b. Efficient
764 grounding DINO: Efficient cross-modality fusion and efficient label assignment
765 for visual grounding in remote sensing. IEEE Transactions on Geoscience and
766 Remote Sensing doi:[10.1109/TGRS.2025.3536015](https://doi.org/10.1109/TGRS.2025.3536015).
- 767 Huang, H., Wang, R., Huang, F., Chen, J., 2025. Analysis and realization
768 of a self-adaptive grasper grasping for non-destructive picking of fruits and
769 vegetables. Computers and Electronics in Agriculture 232, 110119. doi:[10.1016/j.compag.2025.110119](https://doi.org/10.1016/j.compag.2025.110119).
- 770 Jin, T., Han, X., 2024. Robotic arms in precision agriculture: A comprehen-
771 sive review of the technologies, applications, challenges, and future prospects.
772 Computers and Electronics in Agriculture 221, 108938. doi:[10.1016/j.compag.2024.108938](https://doi.org/10.1016/j.compag.2024.108938).
- 773 Jin, Y., Gao, Q., Zhu, X., Yue, Y., Lim, E.G., Chen, Y., Wong, P., Chu,
774 Y., 2024. Vmgnet: A low computational complexity robotic grasping net-
775 work based on vmbamba with multi-scale feature fusion. arXiv preprint
776 arXiv:2411.12520 doi:[10.48550/arXiv.2411.12520](https://doi.org/10.48550/arXiv.2411.12520).

- 780 Jin, Y., Xia, X., Gao, Q., Yue, Y., Lim, E.G., Wong, P., Ding, W., Zhu, X., 2025.
781 Deep learning in produce perception of harvesting robots: A comprehensive
782 review. *Applied Soft Computing* , 112971doi:[10.1016/j.asoc.2025.112971](https://doi.org/10.1016/j.asoc.2025.112971).
- 783 Kebe, G.Y., Higgins, P., Jenkins, P., Darvish, K., Sachdeva, R., Barron,
784 R., Winder, J., Engel, D., Raff, E., Ferraro, F., et al., 2021. A spo-
785 ken language dataset of descriptions for speech-based grounded language
786 learning. *Advances in neural information processing systems* URL: <https://openreview.net/forum?id=Yx9jT3fkBaD>.
- 788 Kumra, S., Joshi, S., Sahin, F., 2020. Antipodal robotic grasping using genera-
789 tive residual convolutional neural network, in: 2020 IEEE/RSJ International
790 Conference on Intelligent Robots and Systems (IROS), IEEE. pp. 9626–9633.
791 doi:[10.1109/IROS45743.2020.9340777](https://doi.org/10.1109/IROS45743.2020.9340777).
- 792 Kuznetsova, A., Rom, H., Alldrin, N., Uijlings, J., Krasin, I., Pont-Tuset, J., Ka-
793 mali, S., Popov, S., Malloci, M., Kolesnikov, A., et al., 2020. The open images
794 dataset v4: Unified image classification, object detection, and visual relation-
795 ship detection at scale. *International journal of computer vision* 128, 1956–
796 1981. URL: <https://storage.googleapis.com/openimages/web/index.html>.
- 798 Li, G., Hammoud, H.A.A.K., Itani, H., Khizbulin, D., Ghanem, B., 2023a.
799 Camel: Communicative agents for "mind" exploration of large scale language
800 model society. doi:[10.48550/arXiv.2303.17760](https://arxiv.org/abs/2303.17760), arXiv:2303.17760.
- 801 Li, S., Yu, H., Ding, W., Liu, H., Ye, L., Xia, C., Wang, X., Zhang, X.P.,
802 2023b. Visual-tactile fusion for transparent object grasping in complex back-
803 grounds. *IEEE Transactions on Robotics* 39, 3838–3856. doi:[10.1109/TR0.2023.3286071](https://doi.org/10.1109/TR0.2023.3286071).
- 805 Li, T., Yan, Y., Yu, C., An, J., Wang, Y., Chen, G., 2024. A comprehensive
806 review of robot intelligent grasping based on tactile perception. *Robotics and*
807 *Computer-Integrated Manufacturing* 90, 102792. doi:[10.1016/j.rcim.2024.102792](https://doi.org/10.1016/j.rcim.2024.102792).
- 809 Li, Z., Liu, J., Li, Z., Dong, Z., Teng, T., Ou, Y., Caldwell, D., Chen, F.,
810 2025. Language-guided dexterous functional grasping by llm generated grasp
811 functionality and synergy for humanoid manipulation. *IEEE Transactions on*
812 *Automation Science and Engineering* doi:[10.1109/TASE.2024.3524426](https://doi.org/10.1109/TASE.2024.3524426).
- 813 Lin, J., Hu, Q., Xia, J., Zhao, L., Du, X., Li, S., Chen, Y., Wang, X., 2023.
814 Non-destructive fruit firmness evaluation using a soft gripper and vision-based
815 tactile sensing. *Computers and Electronics in Agriculture* 214, 108256. doi:[10.1016/j.compag.2023.108256](https://doi.org/10.1016/j.compag.2023.108256).
- 817 Liu, Y., Zhang, J., Lou, Y., Zhang, B., Zhou, J., Chen, J., 2024. Soft bionic
818 gripper with tactile sensing and slip detection for damage-free grasping of
819 fragile fruits and vegetables. *Computers and Electronics in Agriculture* 220,
820 108904. doi:[10.1016/j.compag.2024.108904](https://doi.org/10.1016/j.compag.2024.108904).

- 821 Ma, C., Ying, Y., Xie, L., 2024. Visuo-tactile sensor development and its ap-
822 plication for non-destructive measurement of peach firmness. Computers and
823 Electronics in Agriculture 218, 108709. doi:[10.1016/j.compag.2024.108709](https://doi.org/10.1016/j.compag.2024.108709).
- 824 Rashid, A., Sharma, S., Kim, C.M., Kerr, J., Chen, L.Y., Kanazawa, A.,
825 Goldberg, K., 2023. Language embedded radiance fields for zero-shot task-
826 oriented grasping, in: 7th Annual Conference on Robot Learning. URL:
827 <https://openreview.net/forum?id=k-Fg8JDQmc>.
- 828 Shao, S., Li, Z., Zhang, T., Peng, C., Yu, G., Zhang, X., Li, J., Sun, J., 2019.
829 Objects365: A large-scale, high-quality dataset for object detection, in: Pro-
830 ceedings of the IEEE/CVF international conference on computer vision, pp.
831 8430–8439. URL: www.objects365.org.
- 832 Velasquez, A., Grimm, C., Davidson, J.R., 2024. Compact robotic gripper
833 with tandem actuation for selective fruit harvesting. arXiv preprint
834 arXiv:2408.06674 doi:[10.48550/arXiv.2408.06674](https://doi.org/10.48550/arXiv.2408.06674).
- 835 Visentin, F., Castellini, F., Muradore, R., 2023. A soft, sensorized gripper for
836 delicate harvesting of small fruits. Computers and Electronics in Agriculture
837 213, 108202. doi:[10.1016/j.compag.2023.108202](https://doi.org/10.1016/j.compag.2023.108202).
- 838 Wang, Y.R., Duan, J., Fox, D., Srinivasa, S., 2023. Newton: Are large lan-
839 guage models capable of physical reasoning? arXiv preprint arXiv:2310.07018
840 doi:[10.48550/arXiv.2310.07018](https://doi.org/10.48550/arXiv.2310.07018).
- 841 Xie, W., Valentini, M., Lavering, J., Correll, N., 2024. Deligrasp: Infer-
842 ring object properties with llms for adaptive grasp policies. arXiv preprint
843 arXiv:2403.07832 doi:[10.48550/arXiv.2403.07832](https://doi.org/10.48550/arXiv.2403.07832).
- 844 Xu, J., Jin, S., Lei, Y., Zhang, Y., Zhang, L., 2024. Rt-grasp: Reasoning tuning
845 robotic grasping via multi-modal large language model, in: 2024 IEEE/RSJ
846 International Conference on Intelligent Robots and Systems (IROS), IEEE.
847 pp. 7323–7330. doi:[10.1109/IROS58592.2024.10801718](https://doi.org/10.1109/IROS58592.2024.10801718).
- 848 Xue, L., Gao, M., Xing, C., Martín-Martín, R., Wu, J., Xiong, C., Xu, R.,
849 Niebles, J.C., Savarese, S., 2023. Ulip: Learning a unified representation of
850 language, images, and point clouds for 3d understanding, in: Proceedings of
851 the IEEE/CVF conference on computer vision and pattern recognition, pp.
852 1179–1189. doi:[10.1109/CVPR52729.2023.00120](https://doi.org/10.1109/CVPR52729.2023.00120).
- 853 Yu, X., Tang, L., Rao, Y., Huang, T., Zhou, J., Lu, J., 2022. Point-bert:
854 Pre-training 3d point cloud transformers with masked point modeling, in:
855 Proceedings of the IEEE/CVF conference on computer vision and pattern
856 recognition, pp. 19313–19322. doi:[10.1109/CVPR52688.2022.01871](https://doi.org/10.1109/CVPR52688.2022.01871).
- 857 Zhang, B., Xie, Y., Zhou, J., Wang, K., Zhang, Z., 2020. State-of-the-art
858 robotic grippers, grasping and control strategies, as well as their applications
859 in agricultural robots: A review. Computers and Electronics in Agriculture
860 177, 105694. doi:[10.1016/j.compag.2020.105694](https://doi.org/10.1016/j.compag.2020.105694).

861 Zhang, R., Guo, Z., Zhang, W., Li, K., Miao, X., Cui, B., Qiao, Y., Gao, P.,
862 Li, H., 2022. Pointclip: Point cloud understanding by clip, in: Proceedings of
863 the IEEE/CVF conference on computer vision and pattern recognition, pp.
864 8552–8562. doi:[10.48550/arXiv.2112.02413](https://arxiv.org/abs/2112.02413).

865 Zheng, W., Xie, Y., Zhang, B., Zhou, J., Zhang, J., 2021. Dexterous robotic
866 grasping of delicate fruits aided with a multi-sensory e-glove and manual
867 grasping analysis for damage-free manipulation. Computers and Electronics
868 in Agriculture 190, 106472. doi:[10.1016/j.compag.2021.106472](https://doi.org/10.1016/j.compag.2021.106472).