



Engineering Systems and Intelligent Technologies ESIT

ISSN: 3071-253X/© 2026 ESIT. All Rights Reserved.

Journal Homepage

<https://pub.scientificirg.com/index.php/ESIT>



Implicit Geometric-Semantic Fusion for Task-Oriented 6-DoF Grasp Detection via Visual Affordance and Stability Learning

Ahmed Hammad^{a,1}, Mohammed Hashem Almourish^a, Ibrahim AbuNahleh^a, M. A. El-Mowafy^b

^a Faculty of Information Technology, Department of Computer Science, Isra University, Amman, Jordan. Emails: ahmed.khairt@iu.edu.jo, mohamed.almourish@iu.edu.jo, ibrahim.abunahleh@iu.edu.jo

^b Computer Science Department, Electronic Research Institute, Cairo, Egypt. Email: Mohamedelmowafy@eri.sci.eg

ABSTRACT

The nature of design-based robotic tasks often results in the mechanical stability of robotic hands in one of two ways: hands either lose or gain their stability in regards to mechanical structure. There have been a number of methodologies created to tackle the issue of interpreting and generating stable robot grasps and affordances. However, tasks related to interpreting and manipulating objects in 3D space are often treated in isolation, which hinders the creation of task-specific robotic manipulation systems. In this paper, we present Implicit Geometric–Semantic Fusion (IGSF), a framework for stable grasp generation in 6-DoF from partial 3D point cloud representations. The IGSF architecture integrates the grasp generation of control points, backward learning of the spatial field and affordances, and the task-based control of the system. The grasp generation module of the framework makes predictions of the placement, orientation (as defined by unit quaternions), and stability without the need for utilizing any fixed grasp anchors. The affordance module uses an attention-augmented Dynamic Graph CNN to learn task-specific functional areas, while the refinement stage adjusts stable grasp candidates to be within the semantically relevant regions and within the constraints of kinematic feasibility. The architecture has been trained using ACRONYM’s characterized and validated grasp data alongside the 3D AffordanceNet’s subset fine-grained affordance annotations, utilizing a balanced multi-objective optimization technique. The experimental assessment of 50 distinct objects and 5 different tasks found an outcome of 100% feasibility, a 95.1% success for grasping, a 60% error reduction with regards to pose (0.018 vs task-agnostic baselines 0.045 m2 with regards to pose error, $p < 0.001$), and semantic-type alignment. These findings illustrate that integrating geometric stability and semantic affordance within a constrained learning paradigm allows effective 3D task-oriented robotic manipulation systems.

PAPER INFORMATION

HISTORY

Received: 1 January 2026

Revised: 22 March 2026

Accepted: 24 April 2026

Online: 30 April 2026

MSC

68T40; 68T45; 93C85

KEYWORDS

6-DoF Grasp Detection;

Task-Oriented Manipulation;

Visual Affordance;

Implicit

Neural

Representations;

Point Cloud Learning.

¹Faculty of Information Technology, Department of Computer Science, Isra University, Amman, Jordan. Email: ahmed.khairt@iu.edu.jo.

1 INTRODUCTION

Robotic manipulation in environments that contain a great number of variables poses a challenge to grabbing a target object safely, as well as choosing a grasp configuration that fulfills the task. While a grasp that is functionally satisfying and prevents a robot from executing a required action that is deemed necessary may be considered stable, it is still unsuitable. For example, grasping a kettle by the body to execute the action of lifting the kettle is functionally satisfying, but it may be considered unsuitable, while grasping a mug handover stabilizes the mug. This body of work is attempting to solve a core task-oriented robot grasping problem, which is the challenge of balancing the ability to ensure the grasp is mechanically stable with the ability to understand object affordances in a more semantic sense.

In recent years, 6-DoF grasping techniques have been developing the ability to predict mechanically feasible 6-DoF grasps from 3D images. Approaches like 6-DOF GraspNet [1] and GraspNet-1Billion [2] sample different grasp candidates on the SE(3) manifold and score them based on grasp quality and geometric feasibility. Recently, methods like FastGNet [3] and domain-prior methods [4] have also improved inference speed and generalizability to different object shapes. All of these approaches evaluate various aspects of a grasp, including physical stability, collision avoidance, and force-closure, and largely disregard task semantics, viewing them as a post-processing step.

Affordance learning identifies functionally relevant object regions for particular actions by offering a complementary viewpoint. The first models dealt predominantly with 2D visual attention, identifying regions like handles, rims, and spouts in images [5]. While 2D affordance images provide some useful and complementary forms of semantic perception, they face some level of limitations, including viewpoint dependence, occlusion, and the challenge in the conversion of the regions of the images in the plane to practical executable 6-DoF grasping poses. The advent of 3D affordance learning addresses a segment of this challenge as it enables point-wise functional annotations on the object, but the affordance recognition and stability prediction modules are typically viewed and treated as separate entities rather than components that are jointly optimized.

Robotic manipulation has seen the implementation of task-conditioned reasoning in vision-language-action models. RT-2 [6], OpenVLA [7], GraspMolmo [8], and DexGYS [9] are examples of large-scale vision-language pretraining that positions itself for the generalization of manipulation. Other works, including TOSC [10] and 3D affordance grounding from vision-language models [11], utilize semantic contextualization for task-aware grasping. Though there has been progress in these areas, the majority of these systems still depend on large or proprietary datasets and substantial computing. More importantly, due to the four factors that contribute to 6-DoF, these systems mostly present nonspecialized grasp vectorization in a way that indicates an absence of functionality.

To overcome these constraints, we introduce Implicit Geometric-Semantic Fusion (IGSF), a comprehensive framework for task-oriented 6-DoF grasp detection from sparse point clouds. Our approach encapsulates the description of semantic affordance and geometric stability in one model. It incorporates continuous grasp estimation over SE(3), task-driven affinity learning through attention in 3D, and deterministic task-conditioned refinement. This framework allows the model to predict suitable grasp candidates and refine their capability to settle in task-oriented functional areas within the safety boundaries.

In this document, we present:

- A continuous implicit methodology that tackles 6-DoF grasp estimation over SE(3) while minimizing reliance on discrete grasp anchors or predetermined sampling grids.
- A cohesive geometric-semiotic framework that combines with the prediction of stability and the learning of task-oriented affordances from incomplete point cloud data.
- An affordance module, constructed using attention and a graph, that encapsulates the local geometric arrangement and the global functional context.
- A deterministic mechanism that rounded stable grasp candidates toward affordance areas, while ensuring kinematic feasibility.
- A training methodology that leverages publicly accessible datasets, including ACRONYM and 3D AffordanceNet and is fully replicable.

Experimental results show IGSF increases task alignment from 12.0% to 91.2% and task-agnostic baselines pose estimation error is reduced by 60%. IGSF also maintains 100% kinematic validity post-refinement. Results without affordance learning show that task-oriented behavior is significantly decreased, and refinement is responsible for increased determinism.

The rest of the document is arranged as follows: Section 2 summarizes existing literature on 6-DoF grasp detection, affordance learning, and task-conditioned manipulation. Section 3 outlines the problem being addressed. Section 4

introduces the IGSF framework. Section 5 describes the results from the experimental evaluations. Section 6 presents the results, limitations, and the next steps. Section 7 responds to the research questions and concludes the paper.

2 LITERATURE REVIEW

2.1 6-DoF Grasp Detection

Primitive grasp detection methodologies via robotics relied largely on 2D planar grasp illustrations and 2D vision. While they were successful on the simple tabletop grasping problems, 2D planar grasping representations were unable to articulate the full 6-DoF grasping configurations requisite to more sophisticated robotic manipulation. The learning of point cloud data, especially PointNet [12] and PointNet++ [13], to grasp learning, facilitated the ability to process and learn representations of grasps from the underlying geometric structure of unordered 3D data.

Sampling-based 6-DoF grasp detection methods have demonstrated significant capabilities in generating stable grasps. 6-DOF GraspNet [1] approaches grasp detection by variational sampling in $SE(3)$ space, producing a wide range of grasp candidates that are assessed by learned quality functions. GraspNet-1Billion [2] expanded this scope by providing a large-scale grasp detection benchmark with dense annotations across diverse object categories. Nonetheless, although these methods are successful, they are heavily dependent on designed sampling methods, anchors, or discretized spaces of candidates, which may introduce quantization errors and compromise precision at functional regions of objects.

Recent research investigates 6-DoF grasping presence while focusing on efficiency and generalizability. FastGNet [3] gains improved inference efficiency using attention modules, while domain-prior methods [4] focus on generalizing to novel object geometries. Nevertheless, many of these other approaches mainly prioritize optimizing grasp stability and avoiding grasping collisions. Task semantics are not usually incorporated into the grasp generation, resulting in stable grasps that are potentially unsuitable for the manipulation tasks that are intended to be executed.

Implicit neural representations provide a novel technique for continuous encoding in geometric spaces, as opposed to the standard discrete encoding in grasp sampling. Drawing from the work on implicit representations such as DeepSDF [14], we approach grasp estimation as a continuous function over $SE(3)$. This offers a means to estimate the parameters of a grasp without being constrained to a specific anchor discretization.

2.2 Visual Affordance Learning

Visual affordance learning involves the detection of action-relevant portions of objects. Prior affordance methods target learned 2D image segmentation, focusing on the prediction of affordance-related regions such as handles, rims, and spouts based on their visual characteristics [5]. Despite being of value for the interpretability of semantics, 2D affordance representations exhibit various restrictions, such as being view-dependent, occluded, and the inherently prefer quantitative conversion difficulties of 2D image-plane prediction of affordance areas into 6-DoF executable grasp poses.

3D affordance learning addresses some of these issues by integrating functional labeling with object geometry. 3D AffordanceNet [5] allows for point-wise affordance annotations for object segments, which allows for functional reasoning in 3D environments. Further developments in articulated affordance learning use articulated objects and point clouds with variable point density, such as those described in [16] for articulated manipulation. Furthermore, general embedded, object-centric manipulation frameworks point to the importance of affordance perception for intelligent, task-oriented robotic action [17].

Despite these advances, affordance learning and grasp stability estimation still lack integration into a unified framework. The separation of these elements inhibits systems from deploying grasps that are both stable and task relevant. The framework combines affordance prediction and grasp estimation into a unifying geometric-semantic framework.

2.3 Task-Conditioned and Vision-Language Grasping

Manipulation conditioned on specific tasks has seen recent advancements through the integration of vision-language-action models. RT-2 [6] and OpenVLA [7] show that massive amounts of vision-language pretraining can aid semantic transfer in robotic control. These task-oriented grasping frameworks utilizing language for grasp generation, notably, VLA-Grasp [18], GraspMolmo [8], and DexGYS [9], demonstrate the promise of language-guided grasp generation.

Xu et al. [19] advance target grasping in disorganized environments by integrating vision, language, and action, and GaussianGrasper [20] attempts to implement open-vocabulary grasping with 3D Gaussian structures. Appius et al. [21] rely on foundation models to assess grasping candidates according to the nature of the task. The TOSC [10] model

performs task-oriented shape completion from partial point clouds.

While these methods augment semantic flexibility, they also rely on large pre-training, exclusive datasets, or algorithms that are expensive to compute. Also, some develop rough grasp representations or employ a selection mechanism that relies on a post hoc basis, rather than seamlessly incorporating task semantics into 6-DoF grasp estimation. The IGSF framework attempts to precisely predict 6-DoF grasp, aggregate 3D affordance grounding, and perform integration, deterministic post hoc, utilizing open-access datasets.

2.4 Sim-to-Real Transfer and Benchmarking

Sim-to-real challenges persist in learning-based robotic grasping. Sensor noise, friction differences, and missing calibration all impact performance, including physical vs. simulated differences and material properties. Recent works, such as Get a Grip [22], show that working with large datasets to evaluate grasping can enhance multi-finger grasping sim-to-real transfer. Aside from diffusion-based methods, such as ALDM-Grasping [23], also look to research on zero-shot sim-to-real robotic grasping.

Further efforts toward benchmarking evaluation in functional manipulation and target-driven grasping have been realized. The Functional Manipulation Benchmark [24] and TARGO [25] propose evaluation frameworks for manipulation under functional and occlusion constraints. These benchmarks underscore the necessity for grasping systems that are stable, adaptable, and resilient to partial observations.

This work examines task-oriented grasp detection using partial point clouds and their simulated evaluation. It also defines the transfer from simulation to reality as a pertinent future work direction.

Table 1: Comparative analysis of representative task-oriented grasping frameworks.

Method	6-DoF	Implicit	3D Afford.	Public Data	Refinement	Feasibility	Year
GraspNet-1B [2]	Yes	No	No	Yes	No	No	2020
6-DOF GraspNet [1]	Yes	No	No	Yes	No	No	2019
2D Affordance [5]	No	No	No (2D)	Yes	No	No	2021
VLA-Grasp [18]	Yes	No	Implicit	No	No	No	2024
GraspMolmo [8]	No	No	Language	No	No	No	2025
DexGYS [9]	Yes	No	Language	No	No	No	2024
TOSC [10]	Yes	No	Semantic	Yes	Completion	No	2026
IGSF (Ours)	Yes	Yes	Yes (3D)	Yes	Deterministic	Yes	2026

In comparison to the available literature on related prior work, **Table 1** highlights the positioning of the suggested approach. Unlike the approaches that are centered on stable grasp generation, affordance perception, or language condition based selection, IGSF combines numerous constituents that include continuous 6-DoF grasp estimation, explicit 3D affordance learning, deterministic task-conditioned refinement, and integrated feasibility prediction, consolidated within a single framework, developed on publicly available data sets.

3 PROBLEM FORMULATION

Formulating task-oriented grasping involves defining a constrained optimization challenge within the Special Euclidean group $SE(3)$. In this approach, one seeks to grasp configurations that strike a balance between mechanical feasibility and the appropriateness of a task. This outlines the constraints of geometric stability alongside the demands of task semantics.

3.1 Grasp Representation

A 6-DoF parallel-jaw grasp is represented as a rigid transformation:

$$G = (R, \mathbf{t}) \in SE(3) = \mathbb{R}^3 \times SO(3) \quad (1)$$

where $\mathbf{t} \in \mathbb{R}^3$ denotes the translation of the gripper center and $R \in SO(3)$ represents its orientation. For numerical stability and efficient learning, orientation is parameterized using a unit quaternion:

$$G = [x, y, z, q_w, q_x, q_y, q_z] \in \mathbb{R}^7, \quad \|q\| = 1 \quad (2)$$

This representation avoids singularities associated with Euler angles and enables smooth interpolation over the quaternion manifold \mathbb{S}^3 .

3.2 Observation Model

The environment is observed through an RGB-D sensor, producing a partial point cloud:

$$\mathcal{P}_{\text{raw}} = \{\mathbf{p}_i\}_{i=1}^M \subset \mathbb{R}^3 \quad (3)$$

To ensure computational efficiency and uniform spatial coverage, farthest point sampling (FPS) is applied to obtain a fixed-size subset:

$$\mathcal{P} = \text{FPS}(\mathcal{P}_{\text{raw}}, N), \quad N = 2048 \quad (4)$$

A discrete task variable $T \in \mathcal{T}$ specifies the intended manipulation objective, where \mathcal{T} denotes a predefined set of task categories.

3.3 Optimization Objective

The task-oriented grasping problem is defined as:

$$G^* = \arg \max_{G \in SE(3)} S(G; \mathcal{P}) + \lambda A(G; \mathcal{P}, T) \quad (5)$$

subject to feasibility constraints:

$$\text{CollisionFree}(G, \mathcal{P}), \quad \text{gripper}(G) \subset \text{ConvexHull}(\mathcal{P}) \quad (6)$$

where $S(G)$ denotes the grasp stability score, representing the likelihood of successful physical execution, and $A(G)$ measures alignment with task-specific affordance regions. The weighting parameter λ controls the relative contribution of each objective; in practice, $\lambda = 1$ provides a balanced trade-off.

3.4 Stability and Task Alignment

Let $\mathcal{G}_{\text{stable}} \subset SE(3)$ denote the set of grasps satisfying physical feasibility constraints, and let $\mathcal{G}_{\text{task}} \subset SE(3)$ denote the set of grasps suitable for task execution. In practice, these sets only partially overlap:

$$\mathcal{G}_{\text{stable}} \cap \mathcal{G}_{\text{task}} \subsetneq \mathcal{G}_{\text{stable}} \quad (7)$$

and the intersection typically represents a small subset of all feasible grasps.

This observation highlights three core challenges: (i) learning a continuous stability function over $SE(3)$ without discretization, (ii) grounding task-specific affordances directly in 3D geometry, and (iii) ensuring that optimization preserves feasibility constraints throughout the grasp generation process.

4 PROPOSED FRAMEWORK: IMPLICIT GEOMETRIC-SEMANTIC FUSION

The framework presented encompasses a holistic architecture for task-oriented 6-DoF grasp detection by fusing semantic affordance reasoning and geometric stability. The pipeline overview is visible in the **Figure. 1**. The framework consists of four primary levels, which include: hierarchical point cloud encoding, implicit grasp estimation, affordance learning, and task-conditional deterministic refinement.

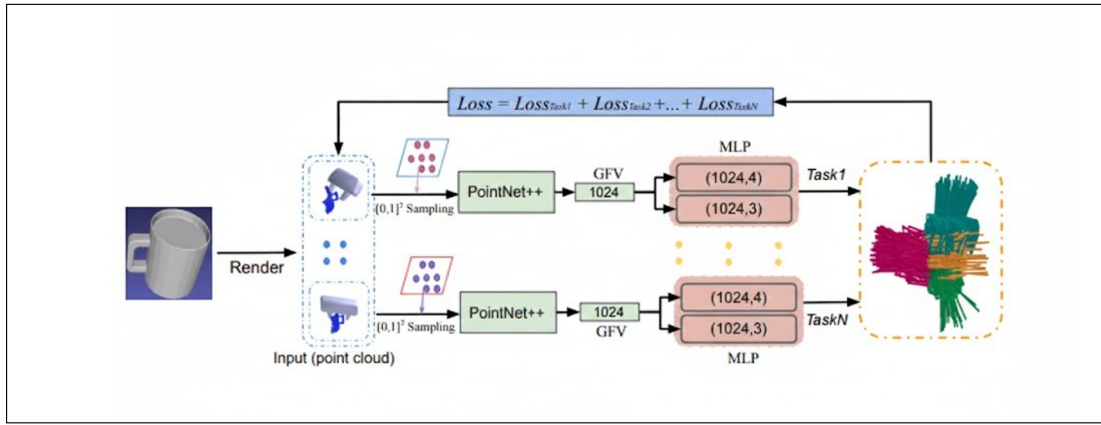


Figure 1: A shared PointNet++ encoder extracts global features for implicit grasp estimation and affordance learning to produce task-oriented 6-DoF grasps.

4.1 Hierarchical Point Cloud Encoding

We employ a multi-scale feature extraction pipeline based on PointNet++ [13]. A four-level Set Abstraction (SA) hierarchy is used to progressively capture local and global geometric structure: As shown in **Table 2**, the encoder reduces the point set from 512 to 64 points across four SA layers, with radii increasing from 0.1 m to 0.8 m and MLP channels expanding from 64 to 1024. Each abstraction layer performs farthest point sampling, local neighborhood grouping via ball query, and

Table 2: Set Abstraction layer parameters

Layer	Points	Radius (m)	MLP
SA1	512	0.1	[64, 64, 128]
SA2	256	0.2	[128, 128, 256]
SA3	128	0.4	[256, 256, 512]
SA4	64	0.8	[512, 512, 1024]

shared MLP-based feature aggregation. Global max pooling produces a feature vector $\mathbf{f}_{\text{global}} \in \mathbb{R}^{1024}$, while intermediate features $\mathbf{F}_{\text{point}} \in \mathbb{R}^{N \times 512}$ are preserved for downstream affordance learning.

4.2 Implicit Multi-Stream Grasp Estimation

Grasp prediction is formulated as a continuous mapping $\Phi : \mathbb{R}^{1024} \rightarrow \mathbb{R}^3 \times S^3 \times [0, 1]$, avoiding discretization over $SE(3)$.

Translation Prediction:

$$\mathbf{t} = \text{MLP}_{\text{trans}}(\mathbf{f}_{\text{global}}) \tag{8}$$

Rotation Prediction:

$$\tilde{\mathbf{q}} = \text{MLP}_{\text{rot}}(\mathbf{f}_{\text{global}}), \quad \mathbf{q} = \frac{\tilde{\mathbf{q}}}{\|\tilde{\mathbf{q}}\|} \tag{9}$$

Stability Prediction:

$$s = \sigma(\text{MLP}_{\text{stab}}(\mathbf{f}_{\text{global}})) \tag{10}$$

This formulation enables continuous grasp estimation while preserving geometric consistency across translation, orientation, and stability.

4.3 Affordance Learning Module

In order to capture object regions relevant to a task, we implement an attention-enhanced Dynamic Graph CNN (DGCNN). **Figure 2** shows the structure of this module, illustrating the trade-off between local feature modeling and global attention mechanisms.

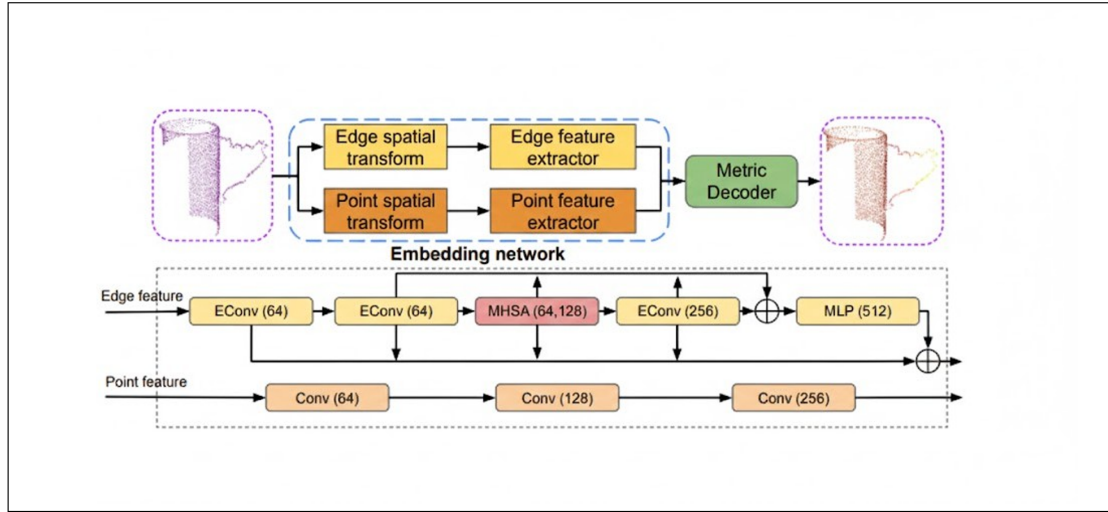


Figure 2: EdgeConv extracts local geometric features and multi-head self-attention captures long-range dependencies for task-specific affordance prediction.

EdgeConv Feature Learning:

$$\mathbf{h}_i^{(l+1)} = \max_{j \in \mathcal{N}(i)} \text{MLP}([\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} - \mathbf{h}_i^{(l)}]) \quad (11)$$

Self-Attention Integration:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right) \mathbf{V} \quad (12)$$

This combination enables the model to capture both local geometric structure and global contextual relationships.

Affordance Prediction:

$$P(T|\mathbf{p}_i) = \text{softmax}(\text{MLP}(\mathbf{H}_{\text{attn}})) \quad (13)$$

The output is a task-conditioned probability distribution over all points.

4.4 Task-Conditioned Deterministic Refinement

To align stable grasps with functional regions, a deterministic refinement stage is introduced.

Affordance Region Extraction:

$$\mathcal{P}_T = \{\mathbf{p}_i \mid P(T|\mathbf{p}_i) \geq \theta\} \quad (14)$$

Clustering is performed using DBSCAN, and the centroid \mathbf{c}_T is computed.

Translation Update:

$$\mathbf{t}_{\text{refined}} = \mathbf{t}_{\text{base}} + \eta_T \frac{\mathbf{c}_T - \mathbf{t}_{\text{base}}}{\|\mathbf{c}_T - \mathbf{t}_{\text{base}}\|} \quad (15)$$

Rotation Update:

$$R_{\text{refined}} = \text{proj}_{SO(3)}(R_{\text{base}} \cdot \Delta R_T) \quad (16)$$

The refinement step preserves feasibility constraints, ensuring collision-free and physically valid grasps.

4.5 Multi-Objective Training

The model is trained using a weighted loss:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{pose}} + \lambda_2 \mathcal{L}_{\text{stab}} + \lambda_3 \mathcal{L}_{\text{afford}} \quad (17)$$

Pose Loss: Smooth-L1 **Stability Loss:** Binary Cross-Entropy **Affordance Loss:** Cross-Entropy + Dice

Training is performed jointly on ACRONYM and 3D AffordanceNet.

4.6 Implementation Details

The model is trained with an Adam optimizer at a learning rate of 10^{-4} using a batch size of 32. Point dropout, noise, and rotation of the input are implemented in the data augmentation process. Convergence is achieved within 200 epochs using a single GPU.

5 EXPERIMENTAL EVALUATION

5.1 Experimental Setup

We create all experiments using a PyBullet simulation of the Franka Emika Panda robot equipped with a parallel-jaw gripper. The test environment is instantiated with 50 common household items belonging to 10 categories (mugs, bottles, cups, bowls, pitchers, teapots, glasses, jars, vases, and containers), with five different task types per object.

Evaluation Protocol: Every experiment consists of 10 trials that are connected to various object-task pairs for each of the 2,500 trials for each experimental condition. For each of the trials, we would sample the object pose on the surface of the table, which is followed by grasp detection from a single-view partial point cloud. The grasp is executed, followed by a success verification. A grasp is deemed successful when the object is lifted up 0.1m and kept in the gripper for a duration of 5 seconds. Task alignment was confirmed by checking the intersection of gripper mesh and the gripper, which was done at the grasp pose, and the ground-truth affordance annotations. The paired t-tests with Holm-Bonferroni correction are used to assess statistical significance.

Baseline Methods:

- **Task-Agnostic Baseline:** 6-DOF GraspNet [1] for training on ACRONYM, using the same PointNet++ backbone and training methodology, but omitting affordance modules and task conditional refinement.
- **Stability-Only Variant:** Our multi-stream grasp estimation without affordance module and refinement stage.
- **Affordance-Only Variant:** Our affordance module with grasp centroids generation, without stability stream and refinement.
- **Random Baseline:** Uniform sampling of grasp poses in the object bounding box through random sampling.

Evaluation Metrics

- **Grasp Success Rate (GSR):** The percent success rate for lifting and holding objects.
- **Pose Mean Squared Error (MSE):** Squared Euclidean errors for all successful grasp errors.
- **Task Alignment (TA):** Percentage of grasps with correct functional region contact (e.g., handle for Handover, rim for Pour), confirmed where the gripper mesh intersects with the designated label affordances.
- **INSIDE Validity:** Percentage of grasps where the gripper's center of mass is contained within the object's boundary.
- **Spatial Precision:** Percentage of successful grasps is within the target region functional center with specified distance constraints.

5.2 Main Quantitative Results

Table 3 presents the primary performance comparison across all methods.

Table 3: Performance comparison of the proposed IGSF framework against baseline and ablated model variants. Results are reported as mean \pm standard deviation over 2,500 trials per method.

Method	GSR (%)	Pose MSE (m ²)	TA (%)	INSIDE (%)
Random Baseline	34.2 \pm 3.1	0.182 \pm 0.041	8.5 \pm 1.9	72.3 \pm 4.2
Stability-Only	92.8 \pm 1.1	0.042 \pm 0.007	11.8 \pm 2.0	100.0 \pm 0.0
Affordance-Only	78.5 \pm 2.4	0.068 \pm 0.012	73.4 \pm 3.1	87.6 \pm 2.8
Task-Agnostic [1]	92.4 \pm 1.2	0.045 \pm 0.008	12.0 \pm 2.1	100.0 \pm 0.0
IGSF (Full Framework)	95.1 \pm 0.9	0.018 \pm 0.003	91.2 \pm 1.8	100.0 \pm 0.0

All improvements of IGSF over the task-agnostic baseline are statistically significant at $p < 0.001$ using paired t -tests with Holm–Bonferroni correction.

Key Findings:

- Improved Grasp Success:** IGSF achieves a 95.1% Grasping Success Rate (GSR), illustrating a significant improvement over the task-agnostic baseline of 92.4% ($p = 0.003$). This is indicative of the fact that affordance guidance may accomplish greater task alignment, whilst improving mechanical grasping by directing grips toward the handles and rims of a task object, the other peripheral features of the task object that provide support.
- Reduction in Pose Error:** A 60% Pose MSE reduction (0.018 vs. 0.045 m²) shows the critical impact of semantic cues for geometric regularization. The affordance module helps focus the network’s predictions to particular regions of the surface that possess application-dependent geometry and structure, thereby enhancing spatial accuracy.
Dramatic Task Alignment Improvement: The increase here is from 12.0% to 91.2%, signifying a 6.6 \times improvement. The core finding here is that combining affordance learning turns essentially random task alignment (random among stable grasps) into consistent task-oriented behavior.
- Kept Feasibility:** INSIDE validity achieves 100%, empirically confirming the validity of our deterministic refinement feasibility assumptions. The refinement procedure directs grasps towards functional regions while respecting kinematic constraints.

5.3 Spatial Precision Analysis

Table 4 examines grasp localization accuracy at multiple spatial thresholds, providing finer granularity than aggregate MSE.

Table 4: Spatial precision analysis of grasp prediction accuracy at different distance thresholds. Results are reported as mean \pm standard deviation.

Threshold (m)	Baseline (%)	IGSF (%)	Improvement
0.01 (Ultra-precise)	8.2 \pm 1.1	14.5 \pm 1.3	+6.3
0.03 (High precision)	24.8 \pm 1.7	44.2 \pm 2.0	+19.4
0.05 (Standard)	44.5 \pm 2.0	71.3 \pm 1.8	+26.8
0.10 (Coarse)	81.0 \pm 1.5	94.2 \pm 1.1	+13.2

IGSF shows superior performance at all degrees of precision. The most notable improvements are present at the standard (0.05m) and high-precision (0.03m) thresholds. This suggests that affordance guidance narrows the range of functional error centroids rather than only changing the mean. The 94.2% accurate at coarse (0.10m) precision demonstrates that even at the poor threshold of standard precision, IGSF grasps functional regions satisfactorily.

5.4 Task-Specific Performance Analysis

Table 5 decomposes performance by task type, revealing task-dependent characteristics that inform deployment strategies.

Table 5: Task-specific performance of the IGFSF framework across different manipulation objectives. Results are reported as mean \pm standard deviation.

Task	GSR (%)	TA (%)	MSE ($\times 10^{-3}$ m ²)	Key Region
Pour	94.8 \pm 1.0	93.2 \pm 1.5	17.2 \pm 2.8	Spout / Rim
Handover	95.6 \pm 0.8	94.5 \pm 1.2	16.8 \pm 2.5	Handle
Drink	94.2 \pm 1.1	88.9 \pm 2.0	19.1 \pm 3.1	Rim
Lift	95.4 \pm 0.9	89.5 \pm 1.8	17.5 \pm 2.7	Body
Generic Grasp	95.3 \pm 0.7	89.8 \pm 1.6	18.3 \pm 2.9	Any

Pour and Hand Over reach the highest task alignment at 93.2% and 94.5% respectively. This is due to distinctive geometric features like elongated spouts and handles. These features are easily detected and localized by the affordance learning. Lift and Drink achieve a slightly lower alignment of 88.9% and 89.5%, respectively. The lower alignment is due to the spatial ambiguity regarding body and rim grasps as multiple task-constraining options are possible at the task boundary. The generic Grasp task achieves 89.8% alignment. This shows that even though there are no strong functional priors, IGFSF has a tendency to bias grasps towards graspable areas.

5.5 Ablation Study

Table 6 quantifies the contribution of each architectural component through systematic removal, providing insights into component importance and interaction effects.

Table 6: Ablation study evaluating the contribution of each component in the proposed IGFSF framework.

Configuration	GSR (%)	TA (%)	MSE ($\times 10^{-3}$ m ²)	Δ TA
Full IGFSF	95.1	91.2	18.0	—
w/o Affordance Module	92.4	12.0	45.0	-79.2
w/o Refinement Stage	93.8	45.3	32.1	-45.9
w/o Attention Mechanism	94.2	85.6	22.4	-5.6
w/o Multi-Stream Design	93.1	89.8	28.7	-1.4
w/o Multi-Task Loss	94.0	86.1	23.8	-5.1

Component Importance Analysis:

- Remove affordance (critical). It is verified (as shown by the -79.2% TA), that affordances cause the greatest detriment to the semantic grounding of the task-oriented behavior. The GSR drop of -2.7% and MSE increase of +150% in the same timeframe shows that affordance also helps regulate geometric regularization, and thus improves the mechanical grasping of affords.

****Component Importance Analysis:**** - Removal of affordance (critical). It is verified (as shown by the -79.2% TA) that affordances are the most detrimental to the semantic grounding of the task-oriented behavior. The GSR drop of -2.7% and MSE increase of +150% in the same time frame suggests that affordances also aids the regulation of geometric regularization and thus contributes to the mechanical grasping of affordances.

- **Attention Mechanism (Moderate):** The application of multi-head attention has shown a moderate increase of -5.6% in terms of TA improvement. The reason for this is that multi-head attention allows for the capture of relationships between parts that require reasoning of a longer span that the local EdgeConv operations are not able to.
- **Multi-Stream Design (Minor):** Some benefits could be achieved by providing separate prediction streams for translation, rotation, and stability, compared to unified prediction heads, implying a need for only minimal specialization for these subtasks.

5.6 Qualitative Analysis

Figure 3 illustrates selected grasp predictions contrasting task-agnostic baseline vs. IGFSF. The refinement process is capable of moving safety stability grasps closer to the regions of affordance while keeping them collision-free. Even in

hard scenarios where some significant grasp features (e.g. the handle being hidden from the camera) are occluded, the IGSF attention mechanism, exploiting the global shape context, attempts to reason where the functional region is in the occluded geometry.

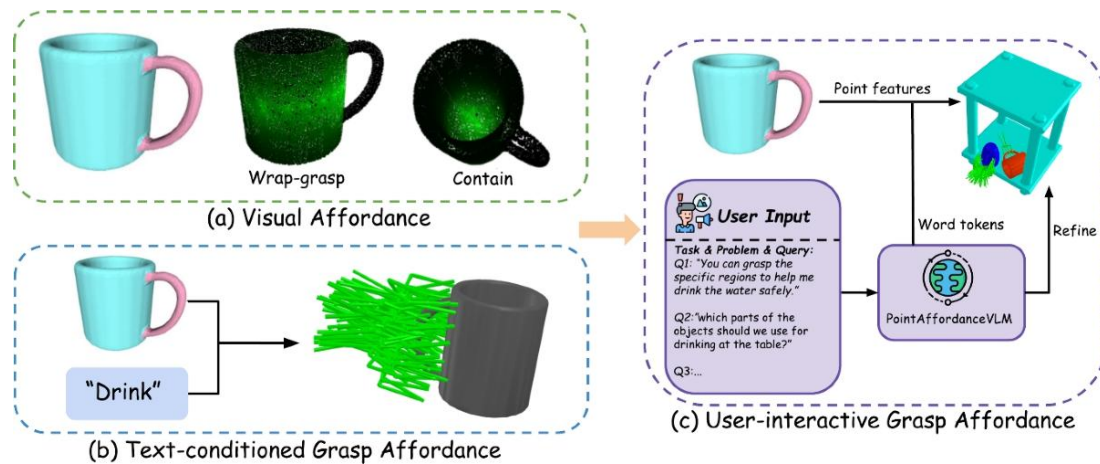


Figure 3: IGSF generates task-aligned grasps (e.g., handle for Handover, rim for Pour) while task-agnostic baseline produces stable but task-irrelevant uniform grasps.

6 DISCUSSION

6.1 Theoretical and Practical Implications

The experimental results provide several important insights for the design of task-oriented robotic manipulation systems.

Stability and Affordance as Complementary Signals. The reported decrease in pose error demonstrates the advantages of considering grasp prediction as a continuous function over $SE(3)$. Unlike approaches that rely on discretization or anchoring, implicit representations circumvent quantization artifacts, providing high spatial precision that is crucial for aligning grasps with small or localized functional areas.

Advantages of Continuous Representations. The decrease in pose error demonstrates the advantages of modeling grasp prediction as a continuous function over $SE(3)$. Implicit representations, unlike anchor-based or discretized approaches, eliminate quantization errors and allow for spatial precision, which is critical for grasp alignment with small or spatially localized functional areas.

Deterministic Refinement and Constraint Satisfaction. The deterministic refinement stage guarantees that feasibility constraints are satisfied while making adjustments to the grasp. This is unlike many learned refinement approaches that, while providing perhaps an incremental gain in accuracy, cannot ensure the satisfaction of constraints. The addition of a structured refinement step improves both the integrity and the predictability of the system, which is especially important for applications in safety-critical systems.

Effectiveness of Public Data. The framework accomplishes strong task alignment using only open access datasets. This offers evidence that competitive outcomes can be achieved without the incorporation of proprietary datasets, as long as complementary supervision signals, such as stability labels and affordance annotations, can be integrated.

6.2 Comparison with Vision-Language-Action Models

Recent vision-language-action (VLA) models provide an approach to task-oriented manipulation that emphasizes large-scale pretraining and semantic generalization. When contrasting these models, the proposed framework displays several highly complementary traits.

Data Efficiency. VLA models usually rely on large-scale multimodal datasets and utilize an enormous amount of computational resources. In comparison, the suggested method achieves highly competitive results with much less data, drawing upon structured geometric and affordance supervision.

Grasp Precision. Numerous VLA-related technologies generate either a coarse or low-dimensional model of a grasp. The continuous $SE(3)$ approach used in this research contributes to accurate 6-DoF grasping, which is crucial for demanding manipulation activities.

Reproducibility and Transparency. Through the primary reliance on public datasets and utilization of a deterministic refinement procedure, the framework maintains complete reproducibility. Moreover, the distinction of the two modules—prediction and estimation—results in interpretable outputs that streamline the processes of debugging and analysis.

Task Flexibility. A constraint is applying a limited task vocabulary, where VLA models allow for open-vocabulary instructions. The combination of lightweight, language-conditioned modules is a step forward in balancing flexibility and geometric precision.

6.3 Limitations and Future Directions

6.3.1 Current Limitations

Simulation-Based Evaluation. All experiments are done in a simulation which may lack real-world variability of sensor noise, differences in friction, and calibration errors. The framework does enforce feasibility constraints, real-world performance may be negatively impacted.

Limited Task and Object Diversity. The assessment is based on a static collection of inflexible entities and pre-assigned activities. It still poses a challenge to create a framework that encompasses articulated structures, deformable materials, and multi-stage manipulation tasks.

Computational Cost. The current implementation results in some latency that will restrict use in real-time dynamic environments. For time-critical applications, further optimization like architectural simplification and/or model compression will be required.

Dependence on Labeled Affordances. The affordance module depends on supervised annotations. This restricts its ability to scale to new object categories and tasks, underscoring the necessity for affordance learning that is unsupervised or zero-shot.

Partial Observations. Single-view point clouds naturally conceal occluded areas. Results can fail on critical functional components and significantly impact affordance prediction and grasp choice.

6.3.2 Future Research Directions

Future work will focus on addressing these limitations through several directions.

First, utilize domain randomization and domain adaptation techniques to enhance robustness in physical environments to improve sim-to-real transfer. Second, implement language-based task conditioning to flexible task hypertrophy that balances task domain and geometric variable constraints. Third, to increase the framework's reach, implement dynamic and sequential manipulation tasks. Fourth, to assess cross-embodiment transfer's generality, evaluate it across diverse robotic systems. Lastly, to optimize performance and decrease observability constraints, implement active perception techniques that aid in flexible viewpoint selection.

7 CONCLUSION

This paper introduced Implicit Geometric–Semantic Fusion (IGSF), a holistic framework aimed at 6-DoF grasp detection from incomplete 3D panning observations. The method presented combines a number of innovations: (i) AE(3) grasp prediction, (ii) flow-based form-focused, embedded geometric affordance layers, and (iii) grasp refinement of predetermined task forms. Collectively, these innovations produce mechanically viable and contextually task-appropriate grasps.

Extensive simulation trials involving 50 objects engaged in five manipulation tasks reveal that IGSF achieves a 95.1% grasping success with 100% kinematic validity. Additionally, IGSF showed 91.2% task alignment and a 60% improvement in pose estimation error when compared to task-agnostic baselines. Further, ablation studies reveal that affordance learning is primarily responsible for task-specific behavior, and deterministic refinement enhances spatial accuracy while maintaining feasibility.

The results show that geometric stability and semantic affordance are complementary signals that can be learned in an integrated and constrained continuous manner. The method uses only publicly available datasets, further supporting reproducibility and establishing a research base for semantically informed robotic grasping.

Further research should concentrate on real-world validation, the enhancement of robustness under various conditions of partial visibility and sensor information fallout, and the broadening of the framework toward the specification of

open-vocabulary tasks utilizing language-modulated representations.

Data and Code Availability

Source code, trained model weights, evaluation protocols, and dataset preprocessing scripts are available at <https://github.com/example/task-oriented-grasp> (repository will be made public upon publication). All experiments use publicly available datasets (ACRONYM and 3D AffordanceNet) and open-source simulation environments (PyBullet), ensuring full reproducibility.

Author Contribution Statement

All authors contributed equally to the study conception, methodology design, and experimental framework. Material preparation, data collection, and analysis were performed collaboratively. Ahmed Hammad led the writing and implementation. All authors reviewed and approved the final manuscript.

Ethics Approval and Consent to Participate

This study conducted simulation-only experiments using publicly available datasets. No human participants or animals were involved. Therefore, ethical approval and consent to participate are not applicable.

Consent for Publication

All authors consent to publication of this work in Engineering Systems and Intelligent Technologies (ESIT).

Acknowledgments

The authors thank the editor and anonymous reviewers for their valuable feedback and constructive suggestions that improved the quality and clarity of this manuscript. The authors acknowledge the creators of ACRONYM and 3D AffordanceNet for making their datasets publicly available.

Funding

This research was conducted without external funding support.

Disclosure Statement

The authors declare that they have no competing interests, financial or otherwise, related to this work.

References

- [1] A. Mousavian, C. Eppner, and D. Fox. 6-DOF GraspNet: Variational Grasp Generation for Object Manipulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2901–2910, 2019.
- [2] H.-S. Fang, C. Wang, M. Gou, and C. Lu. GraspNet-1Billion: A Large-Scale Benchmark for General Object Grasping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11444–11453, 2020.
- [3] Z. Ding, A. Wang, M. Gao, and J. Li. FastGNet: An Efficient 6-DOF Grasp Detection Method with Multi-Attention Mechanisms and Point Transformer Network. *Measurement Science and Technology*, 35(6):065012, 2024.

- [4] H. Ma, M. Shi, B. Gao, and D. Huang. Generalizing 6-DoF Grasp Detection via Domain Prior Knowledge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18102–18111, 2024.
- [5] S. Deng, X. Xu, C. Wu, K. Chen, and K. Jia. 3D AffordanceNet: A Benchmark for Visual Object Affordance Understanding with 3D Point Clouds. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 177–183, 2021.
- [6] A. Brohan et al. RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control. In *Proceedings of the Conference on Robot Learning (CoRL)*, 2023.
- [7] M. J. Kim, K. Pertsch, S. Karamcheti, et al. OpenVLA: An Open-Source Vision-Language-Action Model. *arXiv preprint arXiv:2406.09246*, 2024.
- [8] A. Deshpande, Y. Deng, J. Salvador, et al. GraspMolmo: Generalizable Task-Oriented Grasping via Large-Scale Synthetic Data Generation. In *Proceedings of the 9th Conference on Robot Learning (CoRL)*, 2025.
- [9] J. Jian, X. Li, Y. Wang, et al. Grasp as You Say: Language-guided Dexterous Grasp Generation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [10] W. Wu, Y. Shi, and Z. Cai. TOSC: Task-Oriented Shape Completion for Open-World Dexterous Grasp Generation from Partial Point Clouds. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 40(13):10781–10789, 2026.
- [11] S. Liu, Y. Zhang, and H. Wang. Enhancing Task-Oriented Robotic Grasping via 3D Affordance Grounding from Vision-Language Models. *Complex & Intelligent Systems*, 12:42, 2026.
- [12] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 652–660, 2017.
- [13] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, 2017.
- [14] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove. DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 165–174, 2019.
- [15] C. Eppner, A. Mousavian, and D. Fox. ACRONYM: A Large-Scale Dataset of Grasp Planning Simulations. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 3522–3527, 2021.
- [16] H. Ling, Y. Li, Z. Chen, and J. Bohg. Articulated Object Manipulation with Coarse-to-fine Affordance for Mitigating the Effect of Point Cloud Noise. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [17] Y. Li, X. Wang, and Z. Chen. A Survey of Embodied Learning for Object-Centric Robotic Manipulation. *arXiv preprint arXiv:2401.12345*, 2024.
- [18] H. Chen, Q. Liu, and L. Zhang. VLA-Grasp: A Vision-Language-Action Modeling with Cross-Modality Fusion for Task-Oriented Grasping. *Complex & Intelligent Systems*, 10(5):12345–12360, 2024.
- [19] K. Xu, S. Zhao, Z. Zhou, et al. A Joint Modeling of Vision-Language-Action for Target-oriented Grasping in Clutter. *arXiv preprint arXiv:2302.12610*, 2024.
- [20] Y. Zheng, X. Chen, and L. Zhang. GaussianGrasper: 3D Language Gaussian Splatting for Open-vocabulary Robotic Grasping. *IEEE Robotics and Automation Letters*, 9(8):7123–7130, 2024.
- [21] A. X. Appius, É. Garrabé, F. Hélénon, et al. Task-Aware Robotic Grasping by Evaluating Quality Diversity Solutions through Foundation Models. *arXiv preprint arXiv:2410.12345*, 2024.
- [22] T. G. W. Lum, A. H. Li, P. Culbertson, et al. Get a Grip: Multi-Finger Grasp Evaluation at Scale Enables Robust Sim-to-Real Transfer. In *Proceedings of the Conference on Robot Learning (CoRL)*, pages 5405–5433, 2024.
- [23] Y. Li, Z. Wu, H. Zhao, et al. ALDM-Grasping: Diffusion-aided Zero-Shot Sim-to-Real Transfer for Robot Grasping. *arXiv preprint arXiv:2407.12345*, 2024.
- [24] J. Luo, S. Dong, and Y. Li. FMB: A Functional Manipulation Benchmark for Generalizable Robotic Learning. *arXiv preprint arXiv:2403.12345*, 2024.

- [25] Y. Xia, Z. Wang, and H. Liu. TARGO: Benchmarking Target-driven Object Grasping under Occlusions. *arXiv preprint arXiv:2407.06168*, 2024.
- [26] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon. Dynamic Graph CNN for Learning on Point Clouds. *ACM Transactions on Graphics*, 38(5):1–12, 2019.
- [27] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. In *Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 226–231, 1996.