

MARVO: Marine-Adaptive Radiance-aware Visual Odometry

Anonymous CVPR submission

Paper ID 20972

Abstract

001 Underwater visual localization remains challenging due
002 to wavelength-dependent attenuation, poor texture, and
003 non-Gaussian sensor noise. We introduce MARVO, a
004 physics-aware, learning-integrated Simultaneous Local-
005 ization And Mapping framework that fuses underwater
006 image formation modeling, differentiable matching, and
007 reinforcement-learning optimization. At the front-end, we
008 extend transformer-based feature matcher with a Physics-
009 Aware Radiance Adapter that compensates for color-
010 channel attenuation and contrast loss, yielding geomet-
011 rically consistent feature correspondences under turbid-
012 ity. These semi-dense matches are combined with iner-
013 tial and pressure measurements inside a factor-graph lo-
014 calization backend, where we formulate a keyframe-based
015 visual-inertial-barometric estimator using GTSAM library.
016 Each keyframe introduces (i) Pre-integrated IMU motion
017 factors, (ii) MARVO-derived visual pose factors, and (iii)
018 barometric depth priors, giving a full-state MAP estimate in
019 real time. Lastly, we introduce a Reinforcement-Learning-
020 based Pose-Graph Optimizer that refines global trajec-
021 tories beyond local minima of classical least-squares solvers
022 by learning optimal retraction actions on SE(2). This
023 work highlights the synergy between physics-guided feature
024 adaptation, probabilistic multi-sensor fusion, and learning-
025 based global optimization for robust underwater visual-
026 inertial SLAM.

ther correct the underlying physics of underwater image for-
mation nor effectively fuse uncertain auxiliary sensors such
as pressure and inertial measurements.

We propose the **MARVO**, a *Marine-Adaptive Radiance-Aware Visual Odometry* framework that couples physics-guided front-end perception with probabilistic multi-sensor fusion and offline learning-driven pose-graph refinement. The basic idea of MARVO is that robust underwater VO calls for both (i) perception modules that explicitly compensate for radiometric distortions, and (ii) back-end optimization that can escape the local minima typical of noisy, visually degraded trajectories.

At the front end, MARVO extends LoFTR [30] with a lightweight *Physics-Aware Radiance Adapter* (PARA). PARA modulates intermediate descriptors using learned attenuation and visibility estimates, counteracting color-channel imbalance and restoring feature matchability before transformer attention. This radiance-aware formulation enables stable semi-dense correspondences in regions where standard LoFTR degrades.

These corrected visual factors are fused together with IMU preintegration [11] and barometric depth measures within a keyframe-based factor graph implemented in GT-SAM [8]. Each keyframe contributes with PARA-enhanced visual pose constraints, motion factors and unary depth priors to create a consistent maximum-a-posteriori estimate. The system acts like a fixed-lag smoother, preserving real-time performance while maintaining global consistency across heterogeneous sensing modalities. To further improve long-term stability, MARVO applies a reinforcement-learning (RL) policy in order to refine the SE(2) pose graph; this is inspired by recent work in efficient pose-graph optimization [25]. Instead of solely relying on non-convex least-squares solvers which often get stuck in suboptimal minima, the learned policy proposes retraction actions to steer the trajectory towards globally consistent solutions. This learned refinement complements classical optimization and is particularly effective for turbid or visually sparse scenes.

MARVO integrates physics-informed perception, probabilistic fusion, and learning-based optimization into a reli-

027 1. Introduction

028 Advances in feature matching [30], multi-sensor fu-
029 sion [11] and factor-graph optimization [8] have empow-
030 ered Visual Odometry (VO) and Simultaneous Localiza-
031 tion and Mapping (SLAM) to achieve strong performance in ter-
032 restrial settings. Underwater environments, however, re-
033 main uniquely challenging. Light scattering, wavelength-
034 dependent attenuation and strong non-Gaussian noise pro-
035 duce severe contrast loss, unstable features, and inconsis-
036 tent long-horizon pose estimates. Classical VO and visual-
037 inertial pipelines fail in such conditions because they nei-

038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078

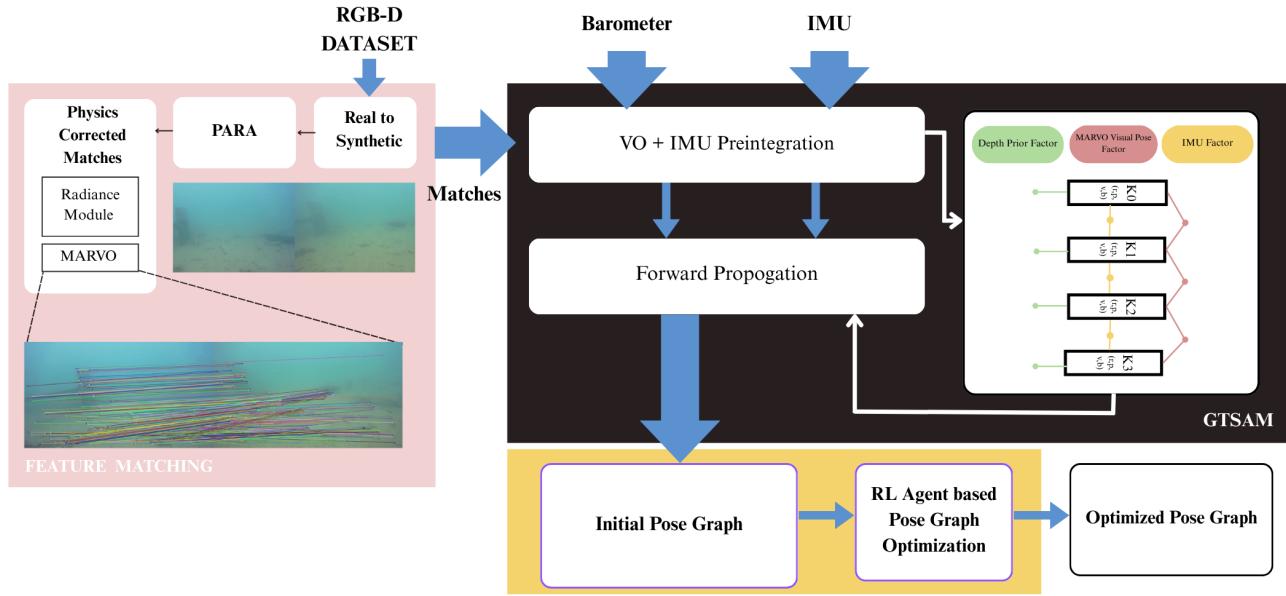


Figure 1. Overview of MARVO. PARA enhances LoFTR features using physically-informed radiance correction. Corrected visual factors are fused with IMU and barometric depth in a GTSAM factor graph to produce real-time VO. An offline reinforcement-learning agent performs pose-graph refinement to obtain globally consistent trajectories.

able underwater localization system. Its contributions are:

- **Physics-aware front end:** a differentiable radiance adapter that compensates wavelength-dependent attenuation within the LoFTR transformer pipeline.
- **Probabilistic multi-sensor fusion:** a visual-inertial-barometric estimator based on a GTSAM fixed-lag smoother with PARA-enhanced constraints.
- **Learning-based global optimization:** an RL-driven pose-graph optimizer that refines SE(2) trajectories beyond the limits of standard least-squares solvers.

Together, these elements provide a single powerful underwater VO system that is able to estimate geometrically consistent trajectories where traditional pipelines fail.

2. Related Works

MARVO leverages three fundamental research directions: feature matching, factor-graph visual-inertial odometry, and reinforcement-learning-based pose-graph optimization.

2.1. Feature Matching

Detector-based methods. Classical local features follow a paradigm of detect-describe-match. While hand-crafted techniques such as SIFT [22], SURF [3], and ORB [29] enjoy robustness to viewpoint and illumination changes, they fail in texturepoor or turbidity degraded underwater imagery. Similarly, learned techniques that improve repeatability via data-driven keypoint detection, such as LIFT [34] and SuperPoint [9], are still limited by the

fundamental sparsity of stable interest points in underwater scenes.

Detector-free methods. Dense matching methods avoid explicit keypoint detection and compute cost volumes or dense correlation fields directly. NCNet [28] and DRC-Net [20] impose neighborhood consensus on dense descriptors, although their convolutional receptive fields limit global reasoning. LoFTR [30] introduced an alternating mechanism of self- and cross-attention in establishing globally consistent semi-dense matches; since then, LoFTR has been widely adopted as a front-end for geometric vision tasks. However, descriptor quality degrades due to contrast loss, backscatter, and wavelength-dependent attenuation underwater, which motivated the physics-aware radiance adaptation in MARVO.

Transformers in geometric vision. Transformers have become foundational in correspondence and motion estimation due to their ability to model long-range relationships. Following ViT [10], attention-based architectures have been deployed in optical flow [31], correspondence estimation [30], and geometric reasoning. MARVO extends this paradigm with a Physics-Aware Radiance Adapter (PARA) that modulates LoFTR’s intermediate features according to learned attenuation cues, restoring discriminability in underwater conditions while retaining global attention benefits.

131 2.2. Factor Graph-Based Visual-Inertial Odometry

132 Pose-graph and factor-graph estimation. Robotic state esti-
 133 mation is often framed as a sparse nonlinear least-squares
 134 problem over poses and landmarks, modeled probabilisti-
 135 cally with factor graphs [8]. Solvers such as GTSAM [8],
 136 g2o [18], and iSAM2 [15] leverage sparsity to enable real-
 137 time inference and form the basis of many state-of-the-art
 138 VO and SLAM methods. Visual-inertial odometry and
 139 preintegration. Combining cameras with IMUs leads to
 140 drift-resistant trajectories, assuming that the system prop-
 141 erly models high-rate inertial signals. IMU preintegration
 142 [11] allows for the efficient incorporation of continuous
 143 IMU measurements in fixed-lag smoothing frameworks
 144 by avoiding full re-integration during optimization. This
 145 has given rise to both highly accurate and computationally
 146 efficient VIO systems [19] suitable for real-time deploy-
 147 ment. MARVO adheres to this paradigm, fusing PARA-
 148 enhanced visual constraints with IMU and barometric depth
 149 in a fixed-lag factor graph. In practical robotic systems,
 150 sensors run at different rates. Continuous-time fusion and
 151 multi-threaded architectures handle asynchronous data ef-
 152 ficiently by interpolating trajectory states and marginaliz-
 153 ing older poses [13, 23]. MARVO uses a similar strat-
 154 egy to enable real-time optimization with heterogeneous vi-
 155 sual-inertial-pressure data.

156 2.3. Reinforcement-Learning-Based Pose-Graph 157 Optimization

158 Classical PGO. Pose-graph optimization (PGO) solves for
 159 globally consistent trajectories by minimizing rotational
 160 and translational constraints over or, typically via Gauss-
 161 Newton or Levenberg-Marquardt with sparse factorizations.
 162 While effective, these solvers remain sensitive to poor ini-
 163 tializations, unreliable loop closures, and strong noise-all
 164 common challenges in underwater environments. RL-based
 165 PGO. Recent work proposes reinforcement-learning-based
 166 optimizers that extend classical solvers by exploring the
 167 pose manifold beyond local gradients. RL-PGO [16] for-
 168 formulates planar PGO as a POMDP and learns retraction ac-
 169 tions on $SE(2)$ that escape suboptimal minima. Distributed
 170 extensions further leverage graph neural networks for multi-
 171 robot settings [17]. MARVO follows a similar philosophy,
 172 using a learned policy to refine pose-graph estimates and
 173 enhance global consistency in the presence of underwater
 174 visual degradation.

175 3. Dataset Augmentation

176 To generate realistic underwater-appearance training data,
 177 we utilise the *SyreaNet* synthesis module [33], which ap-
 178 plies a physics-based underwater image formation model
 179 to in-air RGB-D datasets. We process RGB-depth pairs
 180 from ScanNet [6], TartanAir [32], and Hypersim [27],

181 converting each into a synthetic underwater counterpart
 182 that simulates wavelength-dependent attenuation and depth-
 183 dependent scattering.

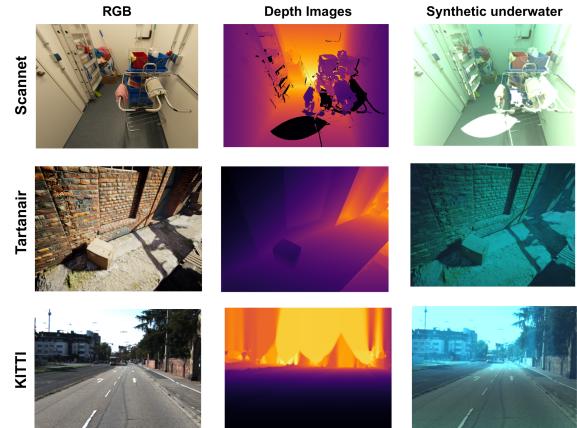


Figure 2. Examples of synthetic data used in MARVO. Each row shows an example from a different dataset (ScanNet, TartanAir, Hypersim): original RGB image, corresponding depth map, and the generated synthetic underwater image with simulated attenuation and scattering.

184 3.1. Physical Model

185 Following the revised underwater image formation
 186 model [1], the observed intensity for color channel $c \in$
 187 $\{R, G, B\}$ is

$$I_c(x) = J_s(x) W_c e^{-\beta_c^D z(x)} + B_c(x), \quad (1) \quad 188$$

189 where $J_s(x)$ is the in-air scene radiance, W_c is the dif-
 190 fuse downwelling term, β_c^D is the attenuation coefficient,
 191 $B_c(x)$ is the backscatter component, and $z(x)$ is the per-
 192 pixel range. This formulation decouples attenuation and
 193 backscatter, enabling physically faithful simulation of haze,
 194 color loss, and contrast degradation under varying water
 195 types.

196 3.2. Synthetic Generation

197 For every RGB-depth pair, SyreaNet samples β_c^D and
 198 $B_c(x)$ from empirical distributions calibrated on real under-
 199 water imagery, emulating different levels of turbidity and
 200 spectral absorption. Illumination variability is emulated by
 201 stochastically perturbing W_c . In total, we synthesize ap-
 202 proximately 120k RGB-D pairs, around 40k per dataset. We
 203 resize all images to 640×480 and normalize them to zero-
 204 mean, unit variance before training. The resulting ren-
 205 derings preserve the geometric structure of the original datasets
 206 but embed realistic radiometric degradation.

207 3.3. Evaluation and Fine-Tuning

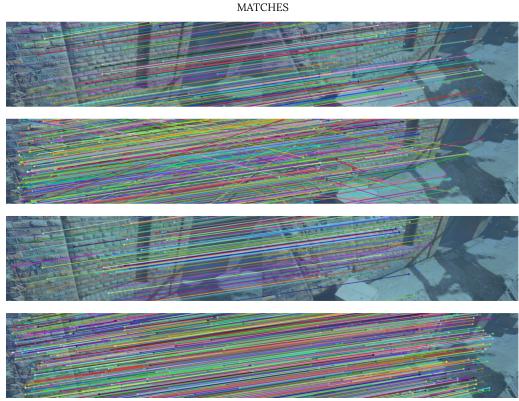


Figure 3. Qualitative feature matching comparison: MARVO produces denser and more geometrically stable correspondences than SuperGlue, ORB, and LoFTR under underwater conditions characterized by turbidity, color attenuation, and low texture. Conventional matchers degrade noticeably, while MARVO maintains semi-dense and spatially coherent matches through physics-aware radiance modulation.

We asses the generalization of the front-end on KITTI [12] in a two-stage protocol: first, perform large-scale pretraining on synthetic underwater data; second, fine-tune on a small subset (10%) of real data. This helps in bridging the synthetic-to-real domain gap while preserving the physical diversity of the synthetic corpus.

We further benchmark the feature-matching pipeline on MegaDepth [21] using the standard image-pair splits and evaluation protocol. MegaDepth provides diverse Internet photo collections with SfM/MVS-derived geometry, allowing us to quantify improvements in match recall and pose accuracy under wide-baseline viewpoint and illumination changes.

4. Feature Matching

At the heart of visual odometry lies feature correspondence. MARVO extends the detector-free matcher LoFTR [30] by adding a physics-guided radiance module inspired from SyreANet [33] for capturing the complex effects of attenuation, scattering, and color imbalance in underwater media. Combining them ensures geometric and photometric consistency even under wavelength-dependent distortions.

4.1. Detector-Free Transformer Backbone

MARVO follows LoFTR’s detector-free paradigm and directly predicts semi-dense correspondences from input image pairs with no explicit keypoint detection. Given rectified images (I_A, I_B), convolutional encoders generate feature maps at 1/8 resolution. The features are enriched with global context across both images by stacked self- and cross-attention layers. High-confidence correspondences

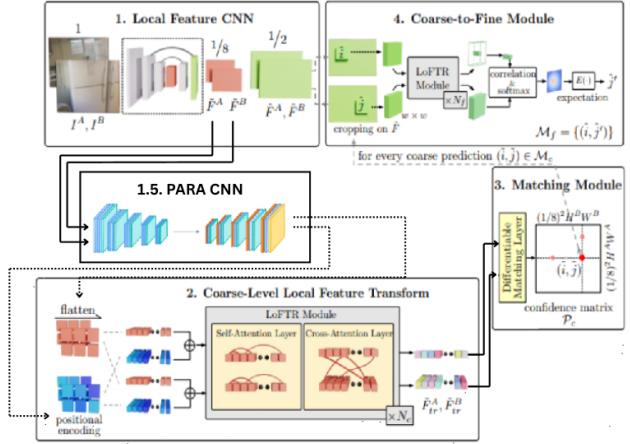


Figure 4. PARA architecture: The Physics-Aware Radiance Adapter takes coarse CNN features and predicts per-pixel attenuation and backscatter fields. These are used to generate a radiance correction mask $\Gamma(x)$, which normalizes intermediate descriptors before transformer matching. PARA compensates for wavelength-dependent attenuation, color imbalance, and contrast degradation common in underwater environments.

are obtained from coarse correlation maps C_{coarse} using dual-softmax and mutual nearest-neighbor filtering and then refined to sub-pixel matches C_{fine} via differentiable correlation. This allows the detector-free architecture to perform robust matching in texture-poor underwater scenarios.

4.2. Physics-Aware Radiance Adapter (PARA)

To adapt transformer matching to underwater imaging, MARVO introduces a *Physics-Aware Radiance Adapter* (PARA) between the CNN encoder and transformer layers.

Let $\mathbf{F} \in \mathbb{R}^{H \times W \times C}$ denote the encoder feature map at 1/8 resolution and let $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$ be the corresponding RGB image, both bilinearly downsampled to the same spatial size. PARA is implemented as a lightweight three-branch prediction head:

$$\mathbf{F}_{\text{shared}} = \phi_{\text{sh}}([\mathbf{F}, \text{BN}(\mathbf{I})]), \quad (2)$$

$$\hat{\beta} \in \mathbb{R}^{H \times W \times 3} = \phi_{\beta}(\mathbf{F}_{\text{shared}}), \quad (3)$$

$$\hat{\mathbf{B}}_{\infty} \in \mathbb{R}^{H \times W \times 3} = \phi_B(\mathbf{F}_{\text{shared}}), \quad (4)$$

$$\hat{\mathbf{z}} \in \mathbb{R}^{H \times W \times 1} = \phi_z(\mathbf{F}_{\text{shared}}), \quad (5)$$

where ϕ_{sh} denotes two 3×3 convolutional layers with ReLU and batch normalization, and $\phi_{\beta}, \phi_B, \phi_z$ are 1×1 convolutional heads that produce per-pixel attenuation, backscatter, and a depth proxy respectively. All predictions are made at the feature resolution, so no extra decoder is required.

We model the underwater image formation for each color channel $c \in \{R, G, B\}$ as

$$I_c(x) = J_c(x) e^{-\beta_c(x)z(x)} + B_{\infty}^c(x)(1 - e^{-\beta_c(x)z(x)}), \quad (6)$$

263 where $J_c(x)$ is the in-air radiance, $\beta_c(x)$ is the attenuation
 264 coefficient, $B_\infty^c(x)$ is the asymptotic backscatter, and
 265 $z(x)$ is the range. During training on synthetic RGB-depth
 266 data, we supervise $\hat{\beta}$ and \hat{B}_∞ using the corresponding Syre-
 267 aNet parameters and provide the ground-truth depth $z(x)$
 268 to PARA. At test time, PARA relies only on the predicted
 269 proxy \hat{z} .

270 To obtain a radiance-corrected estimate of $J_c(x)$ we invert
 271 Eq. (6) using the predicted fields:

$$\hat{J}_c(x) = (I_c(x) - \hat{B}_\infty^c(x)(1 - e^{-\hat{\beta}_c(x)\hat{z}(x)})) \cdot e^{\hat{\beta}_c(x)\hat{z}(x)}. \quad (7)$$

272 From $\hat{J}_c(x)$ we derive a scalar, channel-aggregated correc-
 273 tion mask
 274

$$\Gamma(x) = \frac{1}{3} \sum_{c \in \{R, G, B\}} \frac{\hat{J}_c(x)}{I_c(x) + \epsilon}, \quad (8)$$

275 with a small ϵ to avoid division by zero. The mask $\Gamma(x) \in \mathbb{R}^{H \times W \times 1}$ is then broadcast across channels and applied to
 276 the encoder features as
 277

$$\tilde{\mathbf{F}}(x) = \text{LN}(\Gamma(x) \odot \mathbf{F}(x)), \quad (9)$$

280 where \odot denotes element-wise multiplication and $\text{LN}(\cdot)$ is
 281 layer normalization. The transformer matcher in MARVO
 282 operates on $\tilde{\mathbf{F}}$ instead of \mathbf{F} , so all subsequent attention
 283 layers see features that have been explicitly corrected for
 284 wavelength-dependent attenuation and backscatter.

285 In practice, PARA adds fewer than 5% additional param-
 286 eters relative to the LoFTR backbone and keeps the feature
 287 resolution unchanged, which preserves LoFTR’s computa-
 288 tional profile while significantly improving descriptor con-
 289 sistency in spectrally distorted regions.

290 4.3. Training Objectives

291 The front-end is trained using a combined geometric, pho-
 292 tometric, and physics-based loss:

$$\mathcal{L} = \lambda_{\text{match}} \mathcal{L}_{\text{match}} + \lambda_{\text{photo}} \mathcal{L}_{\text{photo}} + \lambda_{\text{phys}} \mathcal{L}_{\text{phys}}, \quad (10)$$

$$\mathcal{L}_{\text{match}} = \|\hat{\mathbf{P}} - \mathbf{P}^*\|_1, \quad (11)$$

$$\mathcal{L}_{\text{photo}} = 1 - \text{SSIM}(I'_A, I'_B), \quad (12)$$

$$\mathcal{L}_{\text{phys}} = \|\hat{\beta} - \beta_{\text{gt}}\|_1 + \|\hat{B}_\infty - B_{\infty, \text{gt}}\|_1. \quad (13)$$

297 Here β_{gt} and $B_{\infty, \text{gt}}$ are the SyreNet-derived supervision
 298 fields used in Eq. (6). $\mathcal{L}_{\text{match}}$ enforces geometric con-
 299 sistency of correspondences, $\mathcal{L}_{\text{photo}}$ encourages radiance-
 300 corrected view agreement through the PARA-normalized
 301 images I'_A, I'_B , and $\mathcal{L}_{\text{phys}}$ explicitly ties the predicted phys-
 302 ical fields to the underlying underwater image formation
 303 model.

4.4. Synthetic-to-Real Adaptation

MARVO is trained on a combination of synthetic under-
 305 water data and real underwater images. A physics-based
 306 synthesis pipeline generates pairs $(I_{\text{air}}, I_{\text{uw}})$ across multiple
 307 turbidity and illumination conditions. Domain adaptation
 308 and consistency regularization encourage feature invariance
 309 across synthetic-real shifts, thus allowing for robust match-
 310 ing without environment-specific calibration.
 311

4.5. Training Procedure

Our model is trained in two stages: (1) pre-training with
 313 1.2 \times 10⁵ synthetic underwater image pairs from Scan-
 314 Net [6], TartanAir [32] and Hypersim [27]; and (2) fine-
 315 tuning with ~12k real underwater frames including 10%
 316 KITTI [12] and internal field data. The convolutional
 317 layers are partially frozen while fine-tuning transformer
 318 and PARA modules using real-world attenuation statis-
 319 tics. Training utilizes mixed-precision and multi-GPU par-
 320 allelization on 4× NVIDIA A100 GPUs.
 321

The developed physics-aware transformer recovers sta-
 322 ble semi-dense matches even in spectrally distorted and
 323 feature-poor underwater regions, forming the primary con-
 324 straints for downstream GTSAM-based state estimation and
 325 RL-PGO global optimization.
 326

5. Pose Graph Localization

State estimation in MARVO is done by a lightweight factor
 328 graph that fuses: i) physics-aware visual constraints from
 329 PARA–LoFTR, ii) IMU preintegration, and iii) a baromet-
 330 ric depth prior. A fixed-lag smoother is implemented for
 331 the back-end in GTSAM. Unlike typical VIO pipelines,
 332 MARVO involves two new components specific to under-
 333 water applications—a semi-dense visual factor derived from
 334 physics-corrected matches, and a unary depth factor that
 335 suppresses vertical drift.
 336

5.1. State Representation

Every keyframe state \mathbf{x}_i contains orientation, position, ve-
 338 locity, and IMU biases. Only keyframes remain in the optimi-
 339 zation window for real-time operation.
 340

5.2. Graph Factors

IMU Preintegration. We apply standard GTSAM preinte-
 342 gration to provide metric scale and short-term motion con-
 343 straints.
 344

Depth Prior (Underwater). A barometric pressure
 345 reading provides an explicit estimate of depth. We add a
 346 simple unary factor on \mathbf{p}_i an inexpensive but highly effec-
 347 tive constraint that eliminates the vertical drift common in
 348 monocular and low-visibility underwater VO.
 349

MARVO Visual Factor. For every keyframe pair (i, j) ,

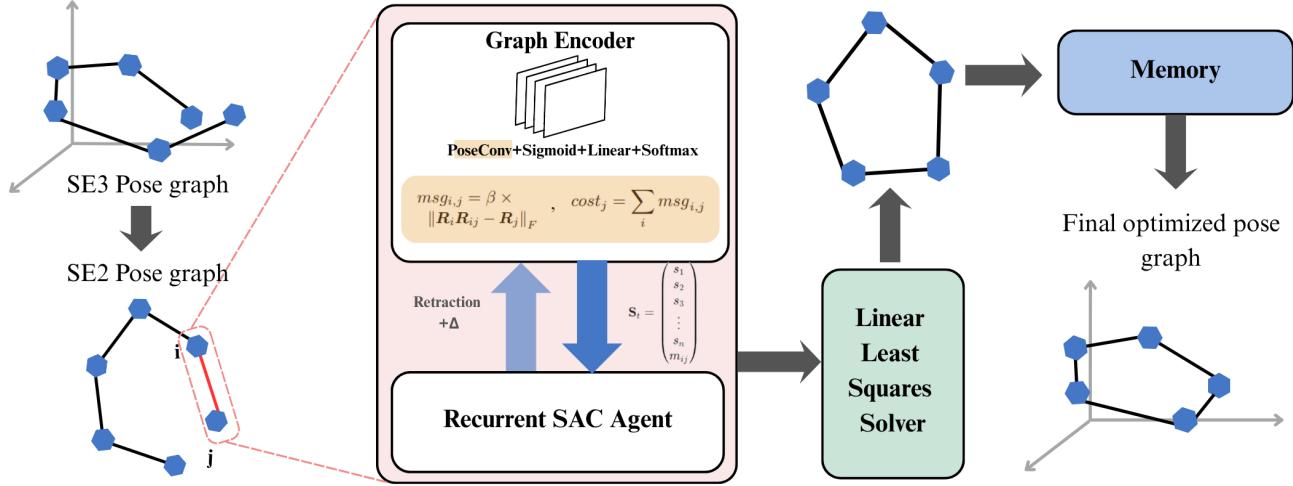


Figure 5. Reinforcement learning-based pose-graph optimization. A GNN encoder maps an initial pose graph to latent edge features, which condition a recurrent SAC agent. The agent iteratively applies retraction actions; a final linear least-squares step produces the optimized pose graph.

351 PARA–LoFTR provides semi-dense matches

$$352 \quad M_{ij} = \{(\mathbf{u}_k^i, \mathbf{u}_k^j)\}_{k=1}^N,$$

353 from which we estimate a relative pose via an essential matrix.
354 A scale variable s_{ij} is co-optimized with the poses,
355 improving robustness under partial stereo loss or poor geo-
356 metric parallax.

357 **Adaptive Covariance.** MARVO weights each visual
358 factor using physics-aware match confidences. The covariance
359 scales inversely with inlier count and spatial coverage,
360 allowing high-visibility frames to dominate while down-
361 weighting degraded or spectrally imbalanced images auto-
362 matically .

363 MARVO’s localization back-end differs from conventional
364 VIO in two aspects: (1) PARA–LoFTR supplies sta-
365 ble, semi-dense constraints even under turbidity and color
366 attenuation, and (2) an underwater-specific depth factor pro-
367 vides drift-free vertical motion. These components produce
368 a reliable initialization for the RL-based global refinement
369 stage that follows.

370 6. Offline Agent-Initiated Pose Graph Optimi- 371 sation

372 6.1. Foundations

373 Classical iterative solvers for pose-graph optimization, such
374 as Gauss–Newton and Levenberg–Marquardt [8, 18], are ef-
375 ficient but prone to local minima under poor initialization or
376 heavy noise. RL-PGO [16] addresses this by casting pose-
377 graph optimization as a partially observable Markov deci-
378 sion process on $\text{SE}(2)$ and learning a policy that applies

379 continuous retraction actions to pose pairs. A message-
380 passing GNN encodes edge residuals [2], enabling the agent
381 to explore beyond local gradients and systematically boot-
382 strap classical solvers on graphs of varying size.

383 6.2. Domain Adaptation for Autonomous Underwa- 384 ter Craft

385 AUVs and ROVs are often roll/pitch stabilized via bal-
386 last and thrusters, leaving yaw as the main rotational
387 degree of freedom, while depth is accurately mea-
388 sured by a pressure sensor [11]. This motivates a
389 dimensionality-reduced formulation: from the GTSAM-
390 based visual–inertial–barometric frontend [7, 15] we extract
391 the planar $\text{SE}(2)$ states (x, y, θ) , build a 2D pose graph on
392 the horizontal plane, and keep roll, pitch, and depth fixed
393 via barometric priors.

394 6.3. Architecture and Integration

395 Our optimizer follows a two-stage pipeline. First, the full
396 $\text{SE}(3)$ factor graph is optimized with a standard solver
397 (iSAM2 or Levenberg–Marquardt) [8, 15] to obtain a con-
398 sistent initial trajectory. We then project this trajectory to
399 $\text{SE}(2)$ and feed the planar graph to the RL agent for refine-
400 ment. After convergence, the corrected (x, y, θ) poses are
401 lifted back to $\text{SE}(3)$ by reattaching the original roll, pitch,
402 and depth estimates, yielding a globally consistent 3D tra-
403 jectory.

404 **Graph Encoding and Message Passing:** We use a
405 GNN encoder that aggregates orientation and translation
406 residuals on each edge $(i, j) \in E$ [2, 16]. Message pass-
407 ing yields per-edge embeddings capturing patterns such as
408 odometric chains and loop closures; these embeddings are
409 pooled into the state fed to the policy.

410 Policy and Sequential Refinement: A recurrent Soft
411 Actor-Critic (SAC) agent with LSTM history [14] selects
412 an edge and outputs a retraction action in $SE(2)$ at each
413 step. Retractions are applied via the exponential map [4],
414 enabling smooth, continuous pose updates while preserving
415 manifold structure.

416 **6.4. Key Innovation: Log-Weighted Orientation** **417 Cost**

418 RL-PGO averages orientation errors uniformly over all
419 edges [16], implicitly assuming equal importance. We in-
420 stead weight each edge’s rotational error by the magnitude
421 of its associated translation, yielding the log-weighted orien-
422 tation cost

$$423 OC_{\log} = \sqrt{\sum_{(i,j) \in E} w_{ij} \|R_i R_{ij} - R_j\|_F^2}, \quad (14)$$

424 with

$$425 w_{ij} = 1 + \beta \log\left(\frac{\|\mathbf{t}_{ij}\|}{\bar{t}} + \epsilon\right), \quad (15)$$

426 where $R_i, R_j, R_{ij} \in SO(2)$ are absolute and relative ro-
427 tations, $\mathbf{t}_{ij} \in \mathbb{R}^2$ is the measured translation, \bar{t} is the mean
428 translation magnitude over all edges, β controls the strength
429 of the weighting, and ϵ is a small constant for numerical sta-
430 bility.

431 The logarithmic form remains monotonically increasing
432 but sublinear in $\|\mathbf{t}_{ij}\|$: long-range constraints are empha-
433 sized without allowing a few very long, noisy edges to
434 dominate the reward. Setting $\beta = 0$ recovers the uniform
435 weighting of [16], while moderate β values better reflect
436 the heterogeneous uncertainty patterns typical in underwa-
437 ter SLAM.

438 **6.5. Post-Optimization Expansion and Final Refine- 439 ment**

440 After RL-based $SE(2)$ refinement, we reattach roll, pitch,
441 and depth from the initial GTSAM estimate (or baromet-
442 ric priors) to obtain a full $SE(3)$ trajectory. A final short
443 Levenberg–Marquardt refinement [18, 25] on the recon-
444 structed 3D pose graph exploits the improved initialization
445 to quickly converge to a high-quality local minimum, even
446 in regimes where purely classical methods tend to stall [5].

447 **7. Experiments**

448 We evaluate MARVO across synthetic underwater bench-
449 marks and real coastal field deployments, comparing
450 against both classical and learning-based VO systems.
451 Our experiments measure (1) correspondence quality un-
452 der wavelength-dependent attenuation and (2) end-to-end
453 odometry accuracy. Due to the absence of standardized
454 underwater VO datasets with ground-truth poses and depth
455 maps, our evaluation incorporates both physically rendered

datasets and real sequences aligned to SfM-based ground
truth. This limitation is inherent to underwater bench-
marking: accurate supervision requires synchronized RGB,
depth, and high-quality pose annotations, which are rarely
available together in public datasets.

461 **7.1. Evaluations**

462 Feature Matching Accuracy. We first benchmark
463 MARVO’s radiance-aware correspondence module against
464 SuperGlue and LoFTR [30]. Using RGB-D datasets ren-
465 dered through a physically based underwater model (per-
466 channel attenuation and scattering), we compute the pose
467 estimation AUC at 5° , 10° , and 20° . MARVO achieves the
468 highest accuracy across all thresholds, showing improved
469 match stability under spectral degradation.

Table 1. Pose estimation AUC on synthetic underwater sequences.

Method	5°	10°	20°
SP + SuperGlue	25.4	42.2	59.7
LoFTR	42.9	59.5	68.2
MARVO (Ours)	49.7	62.9	71.3

470 **Visual Odometry on Synthetic Underwater Dataset.**

471 We next evaluate end-to-end VO performance. All tra-
472 jectories are scale-aligned using a similarity transform.
473 MARVO substantially reduces ATE, angular drift, and rela-
474 tive pose error compared to classical baselines and modern
475 VO pipelines.

Table 2. Synthetic underwater VO performance

Method	ATE (m) \downarrow	RPE (deg/m) \downarrow	Drift (%) \downarrow
ORB-SLAM3	6.45	1.38	5.9
LIBVISO2	5.12	1.14	5.1
MARVO (ours)	2.47	0.61	2.2

Note: Larger-scale multi-sequence evaluations are limited
476 by the scarcity of datasets that provide jointly calibrated
477 ground truth depth, images, and poses, all of which are re-
478 quired to generate physically realistic underwater ren-
479 derings.

481 **7.2. Real-World Field Data**

482 We evaluate MARVO on real underwater imagery collected
483 using a monocular camera, IMU, and depth-pressure sen-
484 sors. Ground truth camera poses are obtained via COLMAP
485 SfM and aligned using a 7-DoF similarity transform. These
486 sequences contain severe turbidity, forward-scattering, and
487 intermittent visibility. MARVO maintains significantly
488 lower drift and angular error compared to classical VO

489 methods, validating its robustness in challenging real-world
490 conditions.

Table 3. VO performance on real underwater field deployments
(Scale Aligned)

Method	ATE (m) ↓	RPE (deg/m) ↓	Drift (%) ↓
ORB-SLAM3	4.12	0.92	3.8
LIBVISO2	3.47	0.85	3.1
MAST3R-SLAM [26]	2.52	0.58	2.2
VGGT-SLAM [24]	2.41	0.56	2.1
MARVO (ours)	1.73	0.34	1.2

491 7.3. Ablation Studies

492 We ablate MARVO to quantify the contribution of each
493 module:

- 494 • **No PARA Module:** Removing physics-guided radiance
495 correction degrades matchability and increases drift, con-
496 firming its necessity under spectral attenuation.
- 497 • **Replace Matcher with LoFTR:** Using vanilla LoFTR
498 instead of PARA-enhanced descriptors produces signifi-
499 cantly worse correspondence reliability.
- 500 • **Replace RL-PGO with Classical PGO:** Gauss–Newton
501 optimization leads to higher drift, especially when loop
502 closures are sparse.
- 503 • **No Physics-Based Radiance Norm:** Using PARA with-
504 out physics-based normalization yields the largest AUC
505 drop, indicating that physics supervision—not merely
506 CNN modulation—is responsible for robustness.

Table 4. Ablation analysis of MARVO components.

Configuration	AUC @10° ↑	ATE (m) ↓	Drift (%) ↓
Full MARVO (ours)	0.92	1.73	1.2
No PARA Module	0.81	2.24	1.9
Replace Feature Matcher (LoFTR)	0.76	2.47	2.3
Replace RL-PGO w/ Classical PGO	0.84	2.05	1.7
No Physics-Based Radiance Norm	0.73	2.68	2.6

507 8. Conclusion

508 We introduced MARVO, a marine-adaptive visual odometry
509 framework that integrates physics-guided feature
510 correction, multi-sensor factor-graph estimation, and
511 reinforcement-learned pose-graph optimization. By em-
512 bedding a wavelength-dependent attenuation and backscat-
513 ter model directly into the transformer correspondence
514 pipeline, MARVO restores descriptor discriminability in vi-
515 sually degraded underwater scenes. Combined with an RL-

enhanced global optimizer and visual–inertial–barometric
frontend, MARVO produces stable, geometrically consistent
trajectories across a wide spectrum of conditions.

Our evaluations demonstrate consistent improvements
over SuperGlue, LoFTR [30], ORB-SLAM3, and LIB-
VISO2 across AUC, ATE, RPE, and drift metrics, both
on synthetic underwater renderings and real-world field de-
ployments. MARVO maintains tracking where classical de-
tectors fail and significantly reduces global trajectory drift
through RL-based refinement.

Limitations and Future Work. A major limitation is
the lack of large-scale underwater VO datasets with ground-
truth trajectories and depth maps. Physics-based rendering
requires metric depth and camera calibration, and accurate
odometry evaluation requires high-quality ground
truth—two properties rarely available simultaneously. This
limits the scope of quantitative benchmarks and neces-
sitates hybrid evaluation using synthetic renderings and
COLMAP-aligned real data.

Future extensions include joint 3D mapping (e.g.,
TSDF fusion or underwater-adapted MVS), learn-
ing full SE(3) global optimization with roll/pitch
coupling, and integrating acoustic depth priors
to handle extreme turbidity or complete visual
dropout.

541 References

- [1] Derya Akkaynak and Tali Treibitz. A revised underwater image formation model. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6723–6732, 2018. 3
- [2] Rana Azzam, Felix H Kong, Tarek Taha, and Yahya Zweiri. Pose-graph neural network classifier for global optimality prediction in 2d slam. *IEEE Access*, 9:80466–80477, 2021. 6
- [3] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. Surf: Speeded up robust features. In *European conference on computer vision*, pages 404–417. Springer, 2006. 2
- [4] José Luis Blanco-Claraco. A tutorial on SE(3) transforma-
tion parameterizations and on-manifold optimization. *arXiv preprint arXiv:2103.15980*, 2021. 7
- [5] Giuseppe Calafiore, Luca Carlone, and Frank Dellaert. Pose graph optimization in the complex domain: Lagrangian duality, conditions for zero duality gap, and optimal solutions. *arXiv preprint arXiv:1505.03437*, 2015. 7
- [6] Angela Dai, Angel X Chang, Manolis Savva, Maciej Hal-
ber, Thomas Funkhouser, and Matthias Nießner. Scannet:
Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017. 3, 5
- [7] Frank Dellaert. Factor graphs and gtsam: A hands-on introduction. *Georgia Institute of Technology, Tech. Rep.*, 2(4), 2012. 6
- [8] Frank Dellaert, Michael Kaess, et al. Factor graphs for robot perception. *Foundations and Trends® in Robotics*, 6(1-2):1–139, 2017. 1, 3, 6

- 571 [9] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabi-
572 novich. Superpoint: Self-supervised interest point detection
573 and description. In *Proceedings of the IEEE conference on*
574 *computer vision and pattern recognition workshops*, pages
575 224–236, 2018. 2
- 576 [10] Alexey Dosovitskiy. An image is worth 16x16 words:
577 Transformers for image recognition at scale. *arXiv preprint*
578 *arXiv:2010.11929*, 2020. 2
- 579 [11] Christian Forster, Luca Carlone, Frank Dellaert, and Da-
580 vide Scaramuzza. On-manifold preintegration for real-time
581 visual–inertial odometry. *IEEE Transactions on Robotics*, 33
582 (1):1–21, 2016. 1, 3, 6
- 583 [12] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel
584 Urtasun. Vision meets robotics: The kitti dataset. *The in-
585 ternational journal of robotics research*, 32(11):1231–1237,
586 2013. 4, 5
- 587 [13] Patrick Geneva, Kevin Eckenhoff, and Guoquan Huang.
588 Asynchronous multi-sensor fusion for 3d mapping and local-
589 ization. In *2018 IEEE international conference on robotics
590 and automation (ICRA)*, pages 5994–5999. IEEE, 2018. 3
- 591 [14] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey
592 Levine. Soft actor-critic: Off-policy maximum entropy deep
593 reinforcement learning with a stochastic actor. In *Inter-
594 national conference on machine learning*, pages 1861–1870.
595 Pmlr, 2018. 7
- 596 [15] Michael Kaess, Hordur Johannsson, Richard Roberts,
597 Viorela Ila, John J Leonard, and Frank Dellaert. isam2: In-
598 cremental smoothing and mapping using the bayes tree. *The
599 International Journal of Robotics Research*, 31(2):216–235,
600 2012. 3, 6
- 601 [16] Nikolaos Kourtzanidis and Sajad Saeedi. RL-pgo: Reinforce-
602 ment learning-based planar pose-graph optimization. *IEEE
603 Control Systems Letters*, 7:3777–3782, 2023. 3, 6, 7
- 604 [17] Sai Krishna Ghanta and Ramviyas Parasuraman. Policies
605 over poses: Reinforcement learning based distributed pose-
606 graph optimization for multi-robot slam. *arXiv e-prints*,
607 pages arXiv–2510, 2025. 3
- 608 [18] Rainer Kümmeler, Giorgio Grisetti, Hauke Strasdat, Kurt
609 Konolige, and Wolfram Burgard. g 2 o: A general frame-
610 work for graph optimization. In *2011 IEEE international
611 conference on robotics and automation*, pages 3607–3613.
612 IEEE, 2011. 3, 6, 7
- 613 [19] Stefan Leutenegger, Simon Lynen, Michael Bosse, Roland
614 Siegwart, and Paul Furgale. Keyframe-based visual–inertial
615 odometry using nonlinear optimization. *The International
616 Journal of Robotics Research*, 34(3):314–334, 2015. 3
- 617 [20] Xinghui Li, Kai Han, Shuda Li, and Victor Prisacariu. Dual-
618 resolution correspondence networks. *Advances in Neural In-
619 formation Processing Systems*, 33:17346–17357, 2020. 2
- 620 [21] Zhengqi Li and Noah Snavely. Megadepth: Learning single-
621 view depth prediction from internet photos. In *Pro-
622 ceedings of the IEEE conference on computer vision and pattern
623 recognition*, pages 2041–2050, 2018. 4
- 624 [22] David G Lowe. Distinctive image features from scale-
625 invariant keypoints. *International journal of computer vi-
626 sion*, 60(2):91–110, 2004. 2
- [23] Jiajun Lv, Xiaolei Lang, Jinhong Xu, Mengmeng Wang,
627 Yong Liu, and Xingxing Zuo. Continuous-time fixed-
628 lag smoothing for lidar-inertial-camera slam. *IEEE/ASME
629 Transactions on Mechatronics*, 28(4):2259–2270, 2023. 3
- [24] Dominic Maggio, Hyungtae Lim, and Luca Carlone. Vggt-
630 slam: Dense rgb slam optimized on the sl (4) manifold. *arXiv
631 preprint arXiv:2505.12549*, 2025. 8
- [25] Gabriel Moreira, Manuel Marques, and Joao Paulo Costeira.
632 Fast pose graph optimization via krylov-schur and cholesky
633 factorization. In *Proceedings of the IEEE/CVF Winter Con-
634 ference on Applications of Computer Vision*, pages 1898–
635 1906, 2021. 1, 7
- [26] Riku Murai, Eric Dexheimer, and Andrew J Davison.
636 Mast3r-slam: Real-time dense slam with 3d reconstruction
637 priors. In *Proceedings of the Computer Vision and Pattern
638 Recognition Conference*, pages 16695–16705, 2025. 8
- [27] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit
639 Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb,
640 and Joshua M Susskind. Hypersim: A photorealistic syn-
641 synthetic dataset for holistic indoor scene understanding. In
642 *Proceedings of the IEEE/CVF international conference on
643 computer vision*, pages 10912–10922, 2021. 3, 5
- [28] Ignacio Rocco, Mircea Cimpoi, Relja Arandjelović, Akihiko
644 Torii, Tomas Pajdla, and Josef Sivic. Neighbourhood con-
645 sensus networks. *Advances in neural information processing
646 systems*, 31, 2018. 2
- [29] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary
647 Bradski. Orb: An efficient alternative to sift or surf. In *2011
648 International conference on computer vision*, pages 2564–
649 2571. Ieee, 2011. 2
- [30] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and
650 Xiaowei Zhou. Loftr: Detector-free local feature matching
651 with transformers. In *Proceedings of the IEEE/CVF con-
652 ference on computer vision and pattern recognition*, pages
653 8922–8931, 2021. 1, 2, 4, 7, 8
- [31] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field
654 transforms for optical flow. In *European conference on com-
655 puter vision*, pages 402–419. Springer, 2020. 2
- [32] Wenshan Wang, Delong Zhu, Xiangwei Wang, Yaoyu Hu,
656 Yuheng Qiu, Chen Wang, Yafei Hu, Ashish Kapoor, and Se-
657 bastian Scherer. Tartanair: A dataset to push the limits of
658 visual slam. In *2020 IEEE/RSJ International Conference
659 on Intelligent Robots and Systems (IROS)*, pages 4909–4916.
660 IEEE, 2020. 3, 5
- [33] Junjie Wen, Jinjiang Cui, Zhenjun Zhao, Ruixin Yan,
661 Zhi Gao, Lihua Dou, and Ben M Chen. Syreanet: A
662 physically guided underwater image enhancement frame-
663 work integrating synthetic and real images. *arXiv preprint
664 arXiv:2302.08269*, 2023. 3, 4
- [34] Kwang Moo Yi, Eduard Trulls, Vincent Lepetit, and Pascal
665 Fua. Lift: Learned invariant feature transform. In *European
666 conference on computer vision*, pages 467–483. Springer,
667 2016. 2