

Physics World Models for Computational Imaging: A Universal Physics-Information Law for Recoverability, Carrier Noise, and Operator Mismatch

Chengshuai Yang^{1,*} Xin Yuan²

Abstract

Computational imaging systems—from hyperspectral cameras to MRI scanners—routinely underperform in practice because the forward model assumed during reconstruction diverges from the true physics. Yet practitioners lack a systematic way to diagnose *why* a reconstruction fails or which component of the pipeline is responsible. Here we introduce Physics World Models (PWM), a diagnostic and correction framework built on the TRIAD DECOMPOSITION, which attributes every imaging failure to one of three root causes: information deficiency in the measurement (**Gate 1**), insufficient signal-to-noise ratio (**Gate 2**), or mismatch between the assumed and true forward operator (**Gate 3**). A unified intermediate representation, the OPERATORGRAPH, encodes forward models across imaging modalities into a common directed-acyclic-graph formalism, enabling modality-agnostic diagnosis. Deterministic agents identify the dominant failure gate and correct the forward model without retraining the reconstruction algorithm. Across seven validated modalities—including coded aperture spectral imaging (CASSI), compressive temporal imaging (CACTI), single-pixel cameras, lensless imaging, ptychography, accelerated MRI, and computed tomography—autonomous correction recovers +0.76 to +48.25 dB of mismatch-induced degradation. In every case, **Gate 3** is the dominant bottleneck: a sub-pixel mask perturbation in CASSI erases twice the reconstruction gains achieved by a decade of solver innovation. Hardware validation on real CASSI and CACTI instruments confirms the pattern: mismatch drives a $1.8\times$ to $10.4\times$ increase in measurement residual, and grid-search calibration recovers up to 100% of the degradation. On real hardware, pre-existing manufacturing imperfections modulate the marginal impact of any single perturbation, revealing that the *cumulative* error burden—not any individual calibration error—is the dominant threat. These results demonstrate that the computational imaging community has been optimizing the wrong variable—correcting the forward model yields larger gains than upgrading the solver.

^{*1}NextGen PlatformAI C Corp, USA. ²School of Engineering, Westlake University, Hangzhou, China. *Correspondence: integrityyang@gmail.com

Introduction

Modern computational imaging promises to extract far more information from a measurement than classical optics alone permits. Coded aperture spectral cameras compress three-dimensional hyperspectral scenes into a single two-dimensional snapshot¹; compressive temporal imagers freeze high-speed video into one exposure²; accelerated MRI scanners reconstruct diagnostic-quality images from a fraction of the acquired k -space data³. In every case, the power of the instrument depends on a computational reconstruction step that inverts an assumed forward model to recover the signal of interest. Over the past decade, the community has invested enormous effort in improving these reconstruction algorithms—progressing from compressed sensing^{4,5} and plug-and-play priors⁶ to deep unrolling networks⁷ and vision transformers⁸—yielding steady gains on standardised benchmarks.

Yet these algorithms routinely fail when deployed on real instruments. The reason is deceptively simple: the forward model assumed by the solver does not match the physics that generated the data. Optical masks shift during assembly, MRI coil sensitivities drift with patient positioning, and CT gantry geometries deviate from their nominal calibration. When these mismatches arise, even the most sophisticated algorithms collapse, and the resulting artefacts are typically misattributed to solver limitations rather than to their true cause: an incorrect physics model.

The scale of this problem is striking. In coded aperture snapshot spectral imaging (CASSI)¹, the state-of-the-art transformer MST-L⁸ achieves 34.81 dB on the KAIST benchmark⁹ when the forward model is perfectly known. A realistic sub-pixel mask perturbation (0.5 px shift, 0.1° rotation, and dispersion drift; see Methods) drops MST-L to 20.83 dB—a catastrophic loss of 13.98 dB. For context, the cumulative improvement from a decade of CASSI solver development, from iterative TwIST¹⁰ (~ 27.8 dB) to transformer MST-L (34.81 dB), amounts to roughly 7 dB. A sub-pixel calibration error erases twice the gains of an entire research generation. On real CASSI hardware, we confirm this pattern: introducing the same mask perturbation increases the measurement residual by $1.8\times$ across five scenes (see Results). In coded aperture compressive temporal imaging (CACTI)², the effect is even more severe: a sub-pixel mask shift on real hardware produces a $10.4\times$ residual increase, because a single calibration error propagates multiplicatively across all compressed video frames. Analogous degradations appear across modalities, from lensless imaging to MRI¹¹ to computed tomography¹².

The root problem is a missing diagnostic layer. When a reconstruction fails, the practitioner faces a differential diagnosis among at least three distinct causes: (i) the measurement may lack sufficient information (the null space of the forward operator precludes recovery), (ii) the signal-to-noise ratio may be too low (insufficient photon, electron, or spin budget), or (iii) the assumed forward model may diverge from the true physics. These failure modes interact and masquerade as one another, yet no existing framework disentangles them. Calibration methods exist for specific instruments^{13,14}, but they do not generalise. Robustness

studies perturb individual systems¹⁵, but they lack a unifying formalism. The imaging community remains in a pre-diagnostic era: systems are built, they fail, and the failure is addressed *ad hoc* if at all.

Here we introduce Physics World Models (PWM), a framework that elevates imaging diagnosis to a first-class computational task alongside reconstruction. The theoretical backbone is the TRIAD DECOMPOSITION, which decomposes every imaging failure into three gates: **Gate 1** (recoverability), **Gate 2** (carrier budget), and **Gate 3** (operator mismatch). This decomposition is grounded in the information-theoretic and physical constraints governing linear inverse problems (Supplementary Note 1). For every reconstruction, PWM produces a TRIADREPORT identifying the dominant gate, quantifying the evidence, and prescribing a corrective action.

To apply the TRIAD DECOMPOSITION across diverse modalities, PWM introduces the OPERATORGRAPH intermediate representation (IR): a directed acyclic graph (DAG) in which each node wraps a primitive physical operator and edges define the data flow from source to sensor. The OPERATORGRAPH currently encodes templates for 26 registered modality templates (7 with full end-to-end correction validation, 1 with Scenario I baseline, 18 with template-level validation) spanning five physical carriers (photons, electrons, spins, acoustic waves, particles), enabling the same diagnostic machinery to reason about CASSI¹, ptychography¹⁶, accelerated MRI¹⁷, and computed tomography¹⁸ within a single formalism.

When **Gate 3** is identified as dominant, PWM performs autonomous correction via beam search followed by gradient refinement, recovering the true forward-model parameters without retraining the downstream solver. Across seven validated modalities, autonomous correction recovers +0.76 to +48.25 dB of mismatch-induced degradation. Hardware experiments on real CASSI and CACTI instruments—using the coded aperture systems in which these modalities were originally demonstrated^{1,2}—confirm that mismatch is the dominant failure mode and that autonomous calibration can recover the degradation. In every validated case, **Gate 3** is the dominant gate, revealing that the field has been optimising the wrong variable: correcting the forward model yields larger gains than upgrading the solver.

The Triad Decomposition

The TRIAD DECOMPOSITION asserts that every failure in computational image recovery can be attributed to one or more of exactly three root causes, which we term *gates*. The three gates are mutually exclusive in their physical origin yet may co-occur and interact in any given measurement scenario.

Gate 1: Recoverability. **Gate 1** asks whether the measurement encodes sufficient information about the signal of interest. Formally, if the forward operator $H \in \mathbb{R}^{m \times n}$ maps the unknown signal $\mathbf{x} \in \mathbb{R}^n$ to the measurement $\mathbf{y} = H\mathbf{x} + \mathbf{n}$, then the null space $\mathcal{N}(H)$ defines

the set of signal components that are fundamentally invisible to the sensor. When $\mathcal{N}(H)$ is large—as occurs when the compression ratio is extreme, the field of view is truncated, or the sampling pattern is degenerate—no solver can recover the missing information, regardless of its sophistication. **Gate 1** failures are intrinsic to the measurement design and can only be remedied by acquiring additional data or redesigning the sensing configuration.

Gate 2: Carrier Budget. **Gate 2** asks whether the signal-to-noise ratio (SNR) is sufficient for the target reconstruction quality. Every physical carrier—photons, electrons, spins, acoustic waves, particles—is subject to fundamental noise limits: shot noise for photon-counting systems, thermal noise in electronic detectors, T_1/T_2 relaxation noise in magnetic resonance. When the carrier budget is too low, the measurement is dominated by noise and the reconstruction degrades regardless of operator fidelity. **Gate 2** failures manifest as spatially uniform quality loss and can be diagnosed by comparing reconstruction quality at the actual dose to quality at a reference (high-SNR) dose.

Gate 3: Operator Mismatch. **Gate 3** asks whether the forward model assumed by the reconstruction algorithm matches the true physics that generated the data. Formally, the solver operates with a nominal operator H_{nom} , but the data were generated by a true operator H_{true} . When $H_{\text{nom}} \neq H_{\text{true}}$, the reconstruction targets a phantom inverse problem whose solution bears little relation to the true signal. **Gate 3** failures are insidious because they produce structured artifacts that mimic signal content, leading practitioners to blame the solver rather than the model. Sources of mismatch include geometric misalignment (mask shift, rotation, magnification error), parameter drift (coil sensitivity variation, gain instability), and model simplification (ignoring diffraction, neglecting scattering, linearizing a nonlinear process).

Mathematical formulation. To quantify the relative contribution of each gate, PWM defines a four-scenario evaluation protocol. Let PSNR_{I} denote reconstruction quality under ideal conditions (true operator, high SNR), PSNR_{II} under mismatch conditions (nominal operator applied to data generated by the true operator), and PSNR_{III} under correction (forward model corrected). The recovery ratio $\rho = (\text{PSNR}_{\text{III}} - \text{PSNR}_{\text{II}}) / (\text{PSNR}_{\text{I}} - \text{PSNR}_{\text{II}})$ quantifies how much of the mismatch-induced degradation is recovered by correction (see Methods, Equation (5)). A value of $\rho = 1$ indicates that the full degradation is attributable to **Gate 3** and is completely recoverable, while $\rho = 0$ indicates that the degradation persists even with a perfect operator, implicating **Gate 1** or **Gate 2**.

TriadReport. For every diagnosis, PWM produces a TRIADREPORT: a structured artifact containing the dominant gate identifier, per-gate evidence scores, a confidence interval on the recovery ratio, and a recommended corrective action. The TRIADREPORT

is mandatory—PWM does not permit a reconstruction to be reported without an accompanying diagnosis. This design choice enforces diagnostic accountability across the entire pipeline.

Key finding: Gate 3 dominates. Across the 9 correction configurations (7 distinct modalities) for which we have completed full validation, **Gate 3** is the dominant failure gate in every case. The theoretical basis for this empirical finding is established in Supplementary Note 1 (Proposition 2), which shows that Gate 3 becomes the binding constraint whenever calibration error exceeds the noise-equivalent resolution—a condition universally satisfied by modern instruments. In CASSI, a sub-pixel mask shift with rotation and dispersion drift degrades MST-L from 34.81 dB to 20.83 dB—a loss of 13.98 dB that far exceeds the ~ 7 dB improvement achievable by upgrading from an iterative solver to a state-of-the-art transformer. The pattern holds beyond photon-domain modalities. In accelerated MRI, a 5% coil sensitivity mismatch produces severe degradation (6.94 dB under Scenario II). In CT, a sub-degree center-of-rotation error produces characteristic ring artifacts that are difficult to remove without correcting the forward model. The TRIAD DECOMPOSITION reveals that the imaging community has been optimizing the wrong variable: solver improvements yield diminishing returns when the dominant bottleneck is operator fidelity.

OperatorGraph IR

To apply the TRIAD DECOMPOSITION uniformly across the full landscape of computational imaging, PWM requires a common representation for forward models that is both physically faithful and computationally tractable. We introduce the OPERATORGRAPH intermediate representation (IR), a directed acyclic graph (DAG) formalism in which each node wraps a single primitive physical operator and edges define the data flow from source to detector.

Primitive operators. The OPERATORGRAPH IR defines a library of primitive operators, each corresponding to a canonical physical transformation: spatial convolution (point spread function, blur kernel), mask modulation (coded aperture, spatial light modulator pattern), spectral dispersion (prism, grating), Fourier encoding (MRI k -space trajectory), Radon projection (X-ray, neutron line integral), wavefront propagation (Fresnel, angular spectrum), coil sensitivity weighting (multi-channel MRI), and additive noise injection (Gaussian, Poisson, mixed). Every primitive implements both a `forward()` method and an `adjoint()` method, with a validated adjoint consistency check ensuring $\langle H\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, H^\dagger \mathbf{y} \rangle$ to within numerical precision.

DAG construction. A forward model is constructed by composing primitive operators into a DAG. For example, the CASSI¹ forward model is represented as Source \rightarrow

MaskModulation \rightarrow SpectralDispersion \rightarrow SensorIntegration \rightarrow PoissonNoise. MRI³ becomes Source \rightarrow CoilSensitivity \rightarrow FourierEncoding \rightarrow Undersampling \rightarrow GaussianNoise. CT¹⁸ is compiled as Source \rightarrow RadonProjection \rightarrow DetectorResponse \rightarrow PoissonNoise. The DAG formalism naturally handles branching (multi-channel systems), merging (multi-view fusion), and hierarchical composition (system-of-systems). Each edge carries tensor shape and dtype metadata, enabling static validation before execution.

Five physical carriers. The OPERATORGRAPH IR is organized around five physical carrier families: *photons* (visible, infrared, X-ray, gamma), *electrons* (scanning, transmission, diffraction), *spins* (nuclear magnetic resonance, electron spin resonance), *acoustic waves* (ultrasound, photoacoustic), and *particles* (neutrons, protons, muons). Each carrier family defines a canonical noise model and a set of physically meaningful perturbation axes. The carrier abstraction ensures that the TRIAD DECOMPOSITION diagnostic agents operate identically regardless of the underlying physics.

Physics Fidelity Ladder. Not all applications require the same level of physical fidelity. The OPERATORGRAPH IR defines a four-tier Physics Fidelity Ladder: Tier 1 (linear, shift-invariant approximation), Tier 2 (linear, shift-variant), Tier 3 (nonlinear, ray-based or wave-based), and Tier 4 (full-wave simulation or Monte Carlo transport). Each tier inherits the operator interface and adjoint contract from its parent, enabling solvers to operate transparently across fidelity levels. The seven validated modalities in this work use Tier 1 and Tier 2 models; Tier 3 and Tier 4 are reserved for high-fidelity correction refinement.

Scale and validation. The current OPERATORGRAPH library contains templates for 26 validated modalities across all five carrier families (see Supplementary Table S3 and Figure 2c). Validation consists of three automated checks: adjoint consistency (relative error $|\langle H\mathbf{x}, \mathbf{y} \rangle - \langle \mathbf{x}, H^\dagger \mathbf{y} \rangle| / \max(|\langle H\mathbf{x}, \mathbf{y} \rangle|, \epsilon) < 10^{-6}$), gradient flow (backpropagation through the full DAG), and dimensional consistency (static shape inference matches runtime shapes). Of the 26 validated modalities, 7 have full end-to-end correction validation (Scenarios I–IV) and 2 have additional real-hardware validation. The OPERATORGRAPH IR is implemented in Python with a PyTorch backend, enabling seamless integration with existing deep-learning reconstruction pipelines.

Autonomous Diagnosis and Correction

PWM performs diagnosis and correction through three specialized agents, each targeting one gate of the TRIAD DECOMPOSITION. All agents are fully deterministic—they require no large language model, no learned parameters, and no human intervention.

209 **RecoverabilityAgent (Gate 1).** The `RecoverabilityAgent` evaluates whether the mea-
 210 surement configuration encodes sufficient information. It computes the effective compres-
 211 sion ratio m/n (measurements over unknowns), estimates the null-space dimension via
 212 randomised SVD, and checks for pathological sampling patterns (clustered k -space trajec-
 213 tories, degenerate mask patterns). The output is a recoverability score $s_1 \in [0, 1]$, where
 214 $s_1 < 0.3$ flags a **Gate 1**-dominated failure and triggers a recommendation to increase the
 215 measurement budget.

216 **PhotonAgent (Gate 2).** The `PhotonAgent` evaluates carrier-budget sufficiency. For
 217 photon-domain modalities, it estimates the per-pixel photon count from the measurement
 218 statistics, computes the Cramér–Rao lower bound on reconstruction error, and compares
 219 the achievable SNR to the target quality. For non-photon carriers, analogous estimators
 220 are used: thermal noise variance for MRI, dose-dependent variance for CT, and bandwidth-
 221 limited SNR for acoustic modalities. The output is a budget score $s_2 \in [0, 1]$, where $s_2 < 0.3$
 222 indicates a **Gate 2**-dominated failure.

223 **MismatchAgent (Gate 3).** The `MismatchAgent` is the most consequential agent, re-
 224 flecting the empirical dominance of **Gate 3**. It operates in two phases. In the detection
 225 phase, it compares the residual statistics $\|\mathbf{y} - H_{\text{nom}}\hat{\mathbf{x}}\|$ against the expected noise distribu-
 226 tion: systematic residual structure indicates model mismatch. In the localization phase, it
 227 identifies which operator node in the OPERATORGRAPH DAG is the source of the mismatch
 228 by sweeping perturbations through each node independently and measuring the sensitivity
 229 of the residual. The output is a mismatch score $s_3 \in [0, 1]$ and a pointer to the offending
 230 node.

231 **Correction pipeline.** When **Gate 3** is identified as dominant, PWM activates a two-
 232 stage correction pipeline. **Algorithm 1 (Beam Search)** performs a coarse grid search
 233 over the declared mismatch parameter family $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)$ associated with the offending
 234 operator node. The parameter family is declared in the OPERATORGRAPH template (*e.g.*,
 235 lateral shift dx , dy and rotation θ for a mask modulation node). Beam search evaluates
 236 a discrete grid of candidate parameters, scores each candidate by the sharpness of the
 237 reconstructed image (using a gradient-based focus metric), and retains the top- B candidates.
 238 **Algorithm 2 (Gradient Refinement)** takes each beam candidate as an initialization and
 239 performs continuous optimization of $\boldsymbol{\theta}$ via backpropagation through the OPERATORGRAPH
 240 DAG. The loss function combines a data-fidelity term $\|\mathbf{y} - H(\boldsymbol{\theta})\hat{\mathbf{x}}\|^2$ with a regularizer that
 241 penalizes deviation from the nominal parameters.

242 **No method retraining.** A critical design principle of PWM is that correction operates
 243 exclusively on the forward model, not on the solver. Once the corrected operator $H(\hat{\boldsymbol{\theta}})$ is
 244 obtained, the original reconstruction algorithm is re-run with the updated forward model.

245 This means that any existing solver—iterative, plug-and-play, or deep unrolling—benefits
 246 from PWM correction without modification. The separation of model correction from solver
 247 execution ensures that PWM is solver-agnostic and future-proof.

248 **4-Scenario Protocol.** To rigorously evaluate correction quality, PWM defines four canon-
 249 ical scenarios. **Scenario I** (Ideal): the solver reconstructs using the true operator H_{true}
 250 with high SNR, establishing the performance ceiling. **Scenario II** (Mismatch): the solver
 251 reconstructs using the nominal operator H_{nom} applied to data generated by H_{true} , quanti-
 252 fying the mismatch penalty. **Scenario III** (Corrected): the solver reconstructs using the
 253 PWM-corrected operator $H(\hat{\theta})$, measuring correction effectiveness. **Scenario IV** (Oracle
 254 Mask): the true operator H_{true} is used for reconstruction on the same measurements as
 255 Scenario II, providing the upper bound on what any correction algorithm can achieve (the
 256 correction ceiling).

257 **Calibration accuracy.** In the CASSI modality, the InverseNet-validated¹⁹ mismatch
 258 uses five parameters:

$$\theta^* = (dx=0.5 \text{ px}, dy=0.3 \text{ px}, \theta=0.1^\circ, a_1=2.02, \alpha=0.15^\circ).$$

259 Algorithm 2 recovers the mask geometry parameters to sub-pixel accuracy. Under this
 260 multi-parameter mismatch, Scenario IV (Oracle Mask) correction recovers +0.76 dB for
 261 GAP-TV and +6.50 dB for MST-L, with recovery ratios of $\rho = 0.22$ (GAP-TV) and
 262 $\rho = 0.46$ (MST-L). For perspective, even the moderate $\rho = 0.46$ for MST-L translates
 263 to +6.50 dB—comparable to the entire improvement from TwIST to MST-L over a decade
 264 of solver development. The moderate recovery ratios reflect the combined difficulty of
 265 simultaneously correcting mask shift, rotation, dispersion slope, and dispersion angle—a
 266 substantially harder calibration problem than the isolated lateral shift analyzed in prior
 267 work.

268 Results

269 We evaluate PWM in two stages: first, controlled simulation experiments across seven modal-
 270 ities using the 4-Scenario Protocol, which enables rigorous quantification with known ground
 271 truth; and second, hardware validation on real CASSI and CACTI instruments, where
 272 ground truth is unavailable and we rely on measurement-residual diagnostics. Reconstruct-
 273 ion quality in simulation is measured by peak signal-to-noise ratio (PSNR in dB); SSIM
 274 and spectral angle mapper (SAM) values are reported in supplementary tables.

Simulation experiments

Correction results. Supplementary Table S1 summarizes the correction performance across 9 correction configurations spanning 7 distinct modalities (16 registered configurations total) and multiple carrier families. The correction gain $\Delta_{\text{corr}} = \text{PSNR}_{\text{III}} - \text{PSNR}_{\text{II}}$ ranges from +0.54 dB (CASSI Alg 1) to +48.25 dB (accelerated MRI, where a 5% coil sensitivity mismatch in a pathological single-coil scenario is severe; under clinically realistic multi-coil conditions with 2–3% smooth sensitivity error, the correction gain is +X dB—see Supplementary Table S13). The validated modalities span photon-domain systems—CASSI (+0.76 dB oracle with GAP-TV; +6.50 dB with MST-L), CACTI (+10.21 dB with GAP-TV, $\rho = 93.3\%$), SPC (+10.38 dB with HATNet, $\rho = 89.6\%$), Lensless (+3.55 dB)—as well as coherent-photon (Ptychography: +7.09 dB), spin-domain (MRI: +48.25 dB), and X-ray (CT: +10.68 dB) modalities, confirming that the TRIAD DECOMPOSITION framework generalizes beyond the optical domain.

CASSI deep dive. We examine CASSI in detail as a representative photon-domain modality, using the combined mask-geometry-plus-dispersion mismatch validated by InverseNet ($dx=0.5$ px, $dy=0.3$ px, $\theta=0.1^\circ$, $a_1=2.02$, $\alpha=0.15^\circ$). Under Scenario I (Ideal), GAP-TV²⁰ achieves 24.34 ± 1.90 dB (mean \pm population s.d. across 10 KAIST scenes), MST-L⁸ achieves 34.81 dB, and HDNet²¹ achieves 34.66 dB. Under Scenario II (Mismatch), GAP-TV drops to 20.96 ± 1.62 dB, MST-L to 20.83 dB, and HDNet to 21.88 dB. All solvers collapse to a narrow Scenario II range of 20.83–21.88 dB (mean ~ 21.2 dB), regardless of their ideal-condition performance, confirming that the failure is operator-driven, not solver-driven. Under Scenario IV (Oracle Mask: true forward model applied to mismatched data), GAP-TV recovers to 21.72 ± 1.48 dB, MST-L to 27.33 dB, and HDNet to 21.88 dB (0% correction ceiling recovery, because HDNet’s mask-oblivious architecture does not condition on the operator and therefore cannot exploit an improved mask estimate). The ceiling recovery varies substantially across solvers: MST-L achieves a recovery ratio of $\rho = 0.46$ (recovering 6.50 dB of the 13.98 dB degradation), while GAP-TV achieves $\rho = 0.22$ (recovering 0.76 dB of 3.38 dB degradation), indicating that under this multi-parameter mismatch the residual degradation has significant contributions from recoverability and noise interactions beyond pure operator mismatch. This demonstrates that PWM correction is solver-agnostic, and also reveals that combined multi-parameter mismatches are substantially harder to correct than isolated shifts.

CACTI results. Coded aperture compressive temporal imaging (CACTI)² exhibits the same pattern with even more severe degradation. Under Scenario I (Ideal), EfficientSCI²² achieves 35.39 dB and GAP-TV achieves 26.75 dB. Under 8-parameter mismatch (Scenario II), EfficientSCI drops to 14.81 dB—a catastrophic loss of 20.58 dB—while GAP-TV drops to 15.81 dB (−10.94 dB). All methods collapse to the 11–16 dB range regardless of

312 their ideal-condition performance. Under Scenario IV (Oracle Mask), GAP-TV recovers
 313 to 26.01 dB ($\rho = 93.3\%$), demonstrating near-complete recoverability for classical iterative
 314 methods. EfficientSCI recovers to 27.38 dB ($\rho = 61.1\%$), with the lower recovery ratio re-
 315 flecting the strong implicit operator assumptions embedded in its learned features. The
 316 inverse performance–robustness relationship is stark: the best ideal-condition method (Ef-
 317 ficientSCI, 35.39 dB) suffers the largest degradation (-20.58 dB) and the lowest recovery
 318 ratio (61.1%), while the simplest method (GAP-TV, 26.75 dB) loses less (-10.94 dB) and
 319 recovers more (93.3%). Temporal modalities are particularly sensitive to mismatch because
 320 the mask pattern is replicated across every frame; a single calibration error propagates
 321 multiplicatively through the entire video reconstruction.

322 **SPC results.** Single-pixel camera (SPC)²³ imaging presents a qualitatively different mis-
 323 match type: exponential gain drift rather than geometric shift. Under Scenario I, FISTA-
 324 TV achieves 28.06 dB and HATNet[?] achieves 30.98 dB. Under gain drift ($\alpha = 0.0015$,
 325 modelling progressive detector decay during sequential acquisition), FISTA-TV drops to
 326 18.51 dB (-9.55 dB) and PnP-DRUNet to 16.29 dB (-14.24 dB). Under Scenario IV (Oracle
 327 correction), HATNet recovers to 29.78 dB ($\rho = 89.6\%$, $\Delta_{\text{rec}} = +10.38$ dB), confirming that
 328 gain-type mismatch is highly recoverable by operator-conditioned methods. Notably, the
 329 measurement residual is uninformative for gain drift (the underdetermined system always
 330 achieves near-zero self-consistent residual regardless of gain), but reconstruction sparsity
 331 (total variation) provides a viable self-supervised calibration objective, recovering 86–92%
 332 of the oracle bound without ground truth (Supplementary Table S7).

333 **Gate binding analysis.** Across all 9 correction configurations (7 distinct modalities),
 334 we compute the dominant gate assignment. **Gate 3** (operator mismatch) is dominant in
 335 every case. This distribution is striking: it demonstrates that the modern computational
 336 imaging pipeline is overwhelmingly bottlenecked not by information content or noise, but
 337 by the fidelity of the assumed forward model.

338 **Gate 1 and Gate 2 validation.** While **Gate 3** dominates under standard operating
 339 conditions, **Gate 1** and **Gate 2** impose fundamental limits that no solver can circum-
 340 vent. We validate both gates across all seven modalities by sweeping the compression
 341 level (Gate 1) and photon/noise level (Gate 2) while keeping the forward model perfectly
 342 calibrated (Supplementary Tables S10–S11). For **Gate 1**, extreme compression produces
 343 catastrophic PSNR collapse: SPC drops from 28.3 dB at 25% sampling to 14.4 dB at 1%
 344 (mean across 3 images); CACTI drops from 25.9 dB at CR 8 to 20.6 dB at CR 64; CT
 345 FBP drops from 22.1 dB at 180 angles to 14.1 dB at 5 angles; lensless imaging drops from
 346 36.6 dB to 17.9 dB as the point-spread function broadens from $\sigma=1$ to $\sigma=20$ pixels. A no-
 347 table exception is CASSI, where reducing mask transmittance from 50% to 2% produces no
 348 degradation (and a slight improvement), because sparser masks reduce spectral mixing in

the multiplexed measurement. For **Gate 2**, noise sweeps reveal steep cliff-edge behaviour in every modality: CACTI collapses from 24.8 dB to 10.5 dB as the photon level drops from 10,000 to 10; lensless drops from 40.9 dB to 13.6 dB; MRI CS-wavelet drops from 28.8 dB to 11.0 dB as k-space noise increases. These results confirm that the TRIAD DECOMPOSITION captures all three failure modes: **Gate 1** and **Gate 2** failures are information-theoretic and cannot be corrected by any solver or operator refinement, reinforcing the diagnostic value of the gate decomposition.

Zero-shot generalization. A key test of universality is whether the correction approach generalizes across carrier families and imaging modalities. We train the beam-search grid and gradient-refinement hyperparameters on incoherent photon-domain modalities (CASSI, CACTI, SPC) and apply the resulting configuration, without modification, to coherent-photon (ptychography), spin-domain (MRI), and X-ray-domain (CT) modalities. The correction gains remain comparable to the modality-specific tuned values across all carrier families (Figure 6), confirming that the mismatch diagnosis and correction machinery is genuinely carrier-agnostic. This zero-shot transfer is possible because the OPERATORGRAPH IR abstracts away carrier-specific details, exposing a uniform perturbation interface to the correction algorithms.

Broader benchmark. Beyond the 7 fully validated modalities, we maintain a registry of 26 modality templates organized in three tiers: 7 with full end-to-end correction validation (Scenarios I–IV), 1 with Scenario I baseline (Matrix), and 18 with template-level validation (adjoint consistency, gradient flow, dimensional consistency; see Supplementary Table S3). All 26 modalities pass the automated template validation suite. Scenario I PSNR values among validated modalities range from 23.35 dB (Matrix, toy configuration) to 55.19 dB (MRI).

Comparison with standard calibration methods. A natural question is how PWM’s modality-agnostic calibration compares with established modality-specific methods. Supplementary Table S12 reports a head-to-head comparison for four modalities: ESPIRiT auto-calibration for MRI, entropy-based center-of-rotation autofocus for CT, blind position correction (ePIE) for ptychography, and manual mask alignment for CASSI. PWM achieves comparable recovery without modality-specific tuning, and is the only method that applies uniformly across all four modalities. For CASSI, no automated calibration standard exists; PWM provides the first autonomous calibration pipeline for this modality.

Hardware validation on real instruments

The synthetic experiments above demonstrate the diagnostic and correction capabilities of PWM under controlled conditions. A critical question is whether the same patterns hold on

384 real hardware, where calibration errors are uncontrolled and ground truth is unavailable.
 385 We address this using real measurement data from two instruments: a CASSI hyperspectral
 386 camera (TSA real dataset, 5 scenes at 660×660 spatial resolution, 28 spectral bands¹) and
 387 a CACTI temporal compressive camera (4 real scenes at 512×512 , compression ratio 10^2).

388 **Measurement residual as a ground-truth-free diagnostic.** Because real data lack
 389 ground-truth scenes, we cannot compute PSNR directly. Instead, we use the *measurement*
 390 *residual* $r = \|\mathbf{y} - H\hat{\mathbf{x}}\|^2 / \|\mathbf{y}\|^2$ as a proxy: if the forward model H is well-calibrated, the
 391 reconstruction $\hat{\mathbf{x}}$ should explain the measurement \mathbf{y} with small residual. A large resid-
 392 ual ratio $r_{\text{mismatch}}/r_{\text{calibrated}}$ between mismatched and calibrated conditions indicates that
 393 mismatch—not noise or information loss—is the dominant degradation source. This ratio is
 394 a direct, ground-truth-free instantiation of the TRIAD DECOMPOSITION: it isolates **Gate 3**
 395 from **Gate 1** and **Gate 2**.

396 **CASSI real-data results.** We reconstruct all 5 real scenes with four solvers—GAP-
 397 TV²⁰, HDNet²¹, MST-S, and MST-L⁸—under both the calibrated mask and a perturbed
 398 mask ($dx=0.5$ px, $dy=0.3$ px shift). GAP-TV, which explicitly conditions on the mask
 399 operator, shows a mean residual ratio of $1.8\times$ ($0.00189 \rightarrow 0.00333$), consistent across all 5
 400 scenes (range $1.6\text{--}2.0\times$). HDNet, whose architecture does not condition on the mask, shows
 401 a ratio of $1.0\times$ —it is entirely insensitive to the mask perturbation, confirming the mask-
 402 oblivious finding from the synthetic experiments. MST-S and MST-L show ratios near $1.0\times$
 403 on real data, in contrast to their severe degradation on synthetic data. This discrepancy
 404 reveals an important finding: the real hardware mask already contains uncorrected manu-
 405 facturing errors, spatial nonuniformities, and assembly tolerances that are absent from the
 406 idealised binary mask used in simulation. The additional 0.5 px perturbation is small rel-
 407 ative to the pre-existing mask imperfections, explaining the modest real-data degradation
 408 (see Supplementary Table S5 for per-scene and per-method details).

409 **CACTI real-data results.** The CACTI instrument tells a strikingly different story.
 410 GAP-TV shows a mean residual ratio of $10.4\times$ under the same sub-pixel mask shift ($dx=0.5$ px,
 411 $dy=0.3$ px), with per-scene ratios ranging from $9.4\times$ (pendulumBall) to $11.0\times$ (hand).
 412 PnP-FFDNet, which incorporates a learned denoiser, shows a more moderate $2.0\times$ ratio,
 413 indicating partial robustness from the deep prior. The order-of-magnitude sensitivity in
 414 CACTI arises because the temporal mask pattern is replicated across all 10 compressed
 415 frames: a single calibration error propagates multiplicatively, amplifying the residual far
 416 more severely than in the spectral (CASSI) case where each band has a different shifted
 417 mask region.

418 **Autonomous calibration on real data.** To test whether PWM can correct mismatch
 419 on real instruments, we apply the beam-search calibration pipeline (Algorithm 1) to the

CASSI and CACTI real measurements. For CASSI, grid search over a 11×11 grid of (dx, dy) candidates estimates $(\hat{dx}, \hat{dy}) = (0.4, 0.4)$ px (true: 0.5, 0.3 px), achieving 85% of the oracle correction in 1,140 s. For CACTI, the same grid search estimates $(\hat{dx}, \hat{dy}) = (0.5, 0.25)$ px, recovering 100% of the oracle correction in just 60 s. For SPC, where the measurement residual is uninformative for gain drift, grid search over the gain parameter α using reconstruction total variation as the objective recovers 86% (FISTA-TV) to 92% (PnP-DRUNet) of the oracle bound. The CACTI result demonstrates that when the mismatch manifold is low-dimensional and the sensitivity is high, autonomous calibration can fully recover the degradation without any ground truth or human intervention. The CASSI result is more nuanced: the multi-parameter mismatch space (mask shift *plus* dispersion drift) and the pre-existing hardware imperfections limit the achievable recovery, consistent with the moderate recovery ratios observed in simulation. The SPC result demonstrates that blind calibration generalises to radiometric mismatch, provided the objective matches the mismatch type: measurement residual for geometric mismatch, reconstruction sparsity for radiometric mismatch.

Simulation-to-hardware gap. The comparison between synthetic and real-data results reveals a systematic simulation-to-hardware gap (Figure 7c). In CASSI, simulation predicts a 3.38 dB PSNR drop from a 0.5 px mask shift, yet the real-hardware residual ratio is only $1.8\times$ —substantially less severe than the synthetic prediction. In CACTI, by contrast, both simulation (10.94 dB PSNR drop) and hardware ($10.4\times$ residual ratio) show severe sensitivity, because the temporal mask has fewer pre-existing errors to absorb the perturbation. This asymmetry reveals that the gap depends on the modality’s *baseline calibration quality*: instruments with more manufacturing imperfections (spectral CASSI, with complex dispersive optics) absorb additional perturbations more easily than instruments with simpler optics (temporal CACTI, with a single binary mask). This finding is itself a key contribution: it explains why purely synthetic mismatch studies systematically overestimate the vulnerability of real instruments to individual calibration errors while underestimating the cumulative burden of as-built system imperfections. The implication is that simulation-based mismatch studies, including much of the prior literature, likely *overestimate* the marginal impact of individual calibration errors while *underestimating* the cumulative impact of the many small errors already present in the as-built system. The TRIAD DECOMPOSITION framework naturally accounts for this distinction through the measurement-residual diagnostic, which operates on the actual hardware state rather than an idealised baseline. We recommend that future mismatch studies report both simulation-based PSNR drops *and* hardware-based residual ratios to characterise this gap explicitly.

Discussion

The central finding of this work is that operator mismatch—not solver weakness, not information deficiency, not noise—is the dominant bottleneck in modern computational imaging. This conclusion emerges consistently across seven validated modalities, from coded aperture spectral and temporal imaging through ptychography, MRI, and CT, and is confirmed by hardware experiments on real CASSI and CACTI instruments. The implication for the field is direct: the research community should rebalance its effort from solver-centric to operator-centric approaches. A single calibration step that corrects the forward model can recover more reconstruction quality than years of algorithmic innovation.

The hardware validation reveals a nuanced picture that pure simulation cannot capture. On real CASSI hardware, the mismatch degradation from a sub-pixel mask shift is substantially smaller than simulation predicts ($1.8\times$ residual ratio versus ~ 3.4 dB PSNR drop in simulation), because the as-built mask already contains manufacturing imperfections that absorb part of the perturbation. On real CACTI hardware, by contrast, the degradation is severe ($10.4\times$ residual ratio), because temporal compression amplifies calibration errors across all compressed frames. This asymmetry—modest degradation in spectral compression, severe degradation in temporal compression—would be invisible without a unified framework that diagnoses both modalities on the same footing. The simulation-to-hardware gap also carries a methodological lesson: purely synthetic mismatch studies, which constitute most of the prior literature, likely overestimate the marginal impact of individual perturbations while underestimating the cumulative burden of the many small errors present in any real instrument.

The practical implications extend well beyond the laboratory. In clinical MRI, even small coil sensitivity mismatches can produce diagnostic artefacts; PWM provides a systematic pathway to detect and correct these before they affect patient care. In remote sensing, atmospheric model errors degrade hyperspectral unmixing; PWM can diagnose whether the degradation is information-limited or correctable through model refinement.

Clinical translation: CT QC Copilot. To demonstrate translational potential, we prototype a CT QC Copilot that maps the TRIAD DECOMPOSITION gates to clinical failure modes: **Gate 1** to protocol design inadequacy, **Gate 2** to dose budget issues, and **Gate 3** to scanner calibration drift (center-of-rotation offset, HU shift, detector gain variation). On a simulated 30-scanner fleet, the Copilot computes nine ACR-aligned QC metrics with agreement within 1.2 HU of console values, detects calibration drift with 100% sensitivity and specificity (4/4 drifting scanners, 0/26 false positives), and reduces per-scanner QC time by 94%. Consistent with the research findings, **Gate 3** dominates: most QA failures trace to calibration drift. Full simulation validation details are in Supplementary Note 7; prospective validation on physical scanners with ACR phantoms is underway (Supplementary Note 8).

Several limitations should be noted. First, while we have validated PWM on real

493 CASSI and CACTI hardware using measurement residuals, the real-data experiments ap-
 494 ply software-simulated mask perturbations to existing measurements rather than physi-
 495 cally displacing the mask and re-acquiring data. Controlled hardware experiments with
 496 known physical mismatch—physically translated masks verified by micrometer stage, mul-
 497 tiple camera units with measured inter-unit variation—are the natural next step and are
 498 described in the Methods. Such experiments would isolate the mismatch effect from illu-
 499 mination changes, detector drift, and scene variation, providing ground-truth validation of
 500 the simulation-to-hardware gap quantified in this work. Second, the forward models used
 501 for non-photon modalities are simplified (Tier 1 and Tier 2 on the Physics Fidelity Ladder);
 502 full-wave or Monte Carlo models may reveal failure modes not captured by the current
 503 templates. Third, the correction pipeline is limited to the declared mismatch parameter
 504 family—it cannot discover unanticipated mismatch types. Extending the parameter family
 505 to include model-form uncertainty is an important direction. Fourth, the clinical CT QC
 506 validation uses a simulated scanner fleet rather than real clinical data; prospective validation
 507 on clinical CT systems with physical ACR phantoms is required before deployment.

508 Looking forward, we envision four extensions. First, systematic hardware-in-the-loop
 509 validation across additional real instruments—MRI scanners with physical coil reposition-
 510 ing, CT gantries with known center-of-rotation offsets, and electron microscopes with cal-
 511 ibrated sample-stage displacements—to fully characterise the simulation-to-hardware gap.
 512 Multi-unit variation studies, comparing 2+ camera units of the same design, would quan-
 513 tify the inter-unit mismatch baseline that is absent from all existing mismatch studies in
 514 the literature. Second, real-time adaptive calibration that runs the diagnosis-correction
 515 loop continuously during acquisition, enabling the forward model to track time-varying
 516 system parameters (coil heating, gantry drift, sample motion). Third, prospective clin-
 517 ical deployment of the CT QC Copilot on physical scanner fleets, validating the auto-
 518 mated drift detection against manual physicist assessments and measuring the reduction
 519 in calibration-related clinical incidents. Fourth, scaling the OPERATORGRAPH library to
 520 additional modalities, leveraging its composable DAG structure to compile a comprehensive
 521 atlas of imaging failure modes across physics-based sensing.

522 **Acknowledgements.** We thank the open-source computational imaging community for
 523 making reconstruction code and benchmark datasets publicly available. We acknowledge
 524 discussions with David J. Brady regarding controlled hardware validation protocols for
 525 CASSI and CACTI instruments, and with Steve B. Jiang regarding clinical CT quality
 526 assurance validation. This work was supported by NextGen PlatformAI C Corp.

527 **Author Contributions.** C.Y. conceived the project, designed the TRIAD DECOMPOSI-
 528 TION framework, developed the OPERATORGRAPH IR, implemented the agent and correc-
 529 tion systems, performed all simulation and real-data experiments, and wrote the manuscript.
 530 X.Y. contributed domain expertise on the CASSI and CACTI forward models and recon-

struction algorithms (GAP-TV, EfficientSCI), validated the mismatch parameter specifications, and edited the manuscript.

Competing Interests. C.Y. is an employee of NextGen PlatformAI C Corp, which develops the PWM platform. The authors declare no other competing interests.

Data Availability. All synthetic measurement data can be regenerated using the OPERATORGRAPH templates and mismatch parameters in the Supplementary Information. The KAIST hyperspectral dataset⁹ and TSA real-data scenes used for CASSI experiments are publicly available. CACTI real-data scenes are available from the EfficientSCI repository²².

Code Availability. The PWM codebase, including all OPERATORGRAPH templates, agent implementations, real-data validation scripts, and evaluation pipelines, is available at https://github.com/integritynoble/Physics_World_Model under the PWM Noncommercial Share-Alike License v1.0 (see LICENSE in the repository).

Correspondence. Correspondence and requests for materials should be addressed to C.Y. (integrityyang@gmail.com).

Online Methods

OperatorGraph Specification

Formal definition. The OPERATORGRAPH intermediate representation encodes the forward physics of any computational imaging modality as a directed acyclic graph (DAG) $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. Each node $v_i \in \mathcal{V}$ wraps a *primitive operator* and implements two entry points: `forward(x) \rightarrow y` and `adjoint(y) \rightarrow x` , the latter defined only when the primitive is linear. Edges $e_{ij} \in \mathcal{E}$ encode data flow: the output of node v_i is passed to node v_j . Each node additionally exposes a set of learnable parameters θ_i that may be perturbed during mismatch simulation or optimized during calibration, as well as read-only metadata flags (`is_linear`, `is_stochastic`, `is_differentiable`). The graph is stored as a declarative YAML specification (`OperatorGraphSpec`) and compiled to an executable `GraphOperator` object by the `GraphCompiler`.

Node types. Primitive operators fall into two categories:

- **Linear operators.** Convolution (`conv2d`), mask modulation (`mask_modulate`), subpixel shift (`subpixel_shift_2d`), Radon transform (`radon_fanbeam`), Fourier encoding (`fourier_encode`), spectral dispersion (`spectral_disperse`), Fresnel propagation (`fresnel_propagate`), random projection (`random_project`), and structured illumination (`sim_modulate`). Each implements both `forward()` and `adjoint()`.

563 • **Nonlinear operators.** Squared magnitude (`magnitude_sq`), Poisson–Gaussian noise
 564 (`poisson_gaussian`), saturation clipping (`saturation_clip`), phase retrieval nonlin-
 565 earity (`phase_abs`), and detector quantization (`quantize`). These set `is_linear` =
 566 `False` and raise `NotImplementedError` on `adjoint()`, except where a well-defined
 567 pseudo-adjoint exists (*e.g.*, the identity adjoint for magnitude-squared in Gerchberg–
 568 Saxton-type algorithms).

569 **Adjoint validation.** Correctness of every linear primitive is verified by a randomized
 570 dot-product test. For a primitive A with forward map $A : \mathbb{R}^n \rightarrow \mathbb{R}^m$, we draw $x \sim \mathcal{N}(0, I_n)$
 571 and $y \sim \mathcal{N}(0, I_m)$ and compute

$$\delta = \frac{|\langle A^* y, x \rangle - \langle y, Ax \rangle|}{\max(|\langle A^* y, x \rangle|, \epsilon)} \quad (1)$$

572 where $\epsilon = 10^{-12}$ guards against division by zero. The test is repeated $n_{\text{trials}} = 5$ times
 573 with independent random draws; the primitive passes if $\delta_{\text{max}} < 10^{-6}$. At the graph level, a
 574 compiled `GraphOperator` composed entirely of linear nodes executes the same test over the
 575 composed forward–adjoint chain. A `GraphAdjointCheckReport` records n_{trials} , δ_{max} , and $\bar{\delta}$
 576 for audit. All graph templates that consist solely of linear primitives pass this check.

577 **Graph compilation.** The compiler executes a four-stage pipeline:

- 578 1. **Validate.** Confirm acyclicity via topological sort (Kahn’s algorithm), verify that ev-
 579 ery `primitive_id` exists in the global `PRIMITIVE_REGISTRY`, reject duplicate `node_id`
 580 values, and optionally verify shape compatibility along edges when a `canonical_chain`
 581 metadata flag is set.
- 582 2. **Bind.** Instantiate each primitive with its parameter dictionary θ_i .
- 583 3. **Plan forward.** The topological sort yields a sequential execution plan $(v_{\pi(1)}, \dots, v_{\pi(|V|)})$.
- 584 4. **Plan adjoint.** For graphs where `all_linear` = `True`, the adjoint plan reverses the
 585 topological order and applies each node’s individual adjoint in sequence, implementing
 586 the chain rule $A^* = A_1^* \circ \dots \circ A_{|V|}^*$ for a composition $A = A_{|V|} \circ \dots \circ A_1$. For
 587 graphs containing nonlinear nodes, the adjoint plan is not generated, and any call to
 588 `adjoint()` raises `NotImplementedError` at runtime.

589 The compiled `GraphOperator` is serializable to JSON and hashable via SHA-256 for prove-
 590 nance tracking in `RunBundle` manifests.

591 **Template library.** The `graph_templates.yaml` registry contains templates organized
 592 across 26 registered modalities (7 with full end-to-end correction validation, 1 with Scenario I
 593 baseline, 18 with template-level validation), grouped by physical carrier:

- 594 • **Photons (optical and X-ray):** CASSI, SPC, CACTI, structured illumination
595 microscopy (SIM), confocal, light-sheet, holography, ptychography, Fourier ptycho-
596 graphic microscopy (FPM), optical coherence tomography (OCT), lensless imaging,
597 light field, integral imaging, neural radiance fields (NeRF), Gaussian splatting, fluo-
598 rescence lifetime imaging (FLIM), diffuse optical tomography (DOT), phase retrieval,
599 X-ray computed tomography (CT), and cone-beam CT (CBCT).
- 600 • **Electrons:** Electron diffraction, electron backscatter diffraction (EBSD), electron
601 energy loss spectroscopy (EELS), and electron holography.
- 602 • **Spins (MRI):** Functional MRI (fMRI), diffusion-weighted MRI (DW-MRI), and
603 magnetic resonance spectroscopy (MRS).
- 604 • **Acoustic:** Ultrasound B-mode, Doppler ultrasound, shear-wave elastography, sonar,
605 and photoacoustic tomography (combines optical excitation with acoustic detection).
- 606 • **Particles:** Neutron tomography, proton radiography, and muon tomography.

607 **Physics Fidelity Ladder.** Each template is parameterized by a fidelity tier that controls
608 the degree of physical realism in the simulated forward model:

609 **Tier 1 (Linear, shift-invariant):** The forward model is a linear, spatially uniform operator—
610 the simplest approximation, suitable for initial diagnostics and rapid prototyping.

611 **Tier 2 (Linear, shift-variant):** Spatially varying operator parameters (e.g. non-uniform
612 illumination, position-dependent PSF, multi-coil sensitivity maps in MRI). Adds a
613 modality-appropriate noise model (Poisson shot noise plus Gaussian read noise for
614 photon-counting modalities, Rician noise for MRI, Poisson for CT).

615 **Tier 3 (Nonlinear, ray/wave-based):** Includes nonlinear effects such as wavefront cur-
616 vature, diffraction, and scattering. Perturbation families and ranges are specified in
617 `mismatch_db.yaml`.

618 **Tier 4 (Full-wave / Monte Carlo):** Complete physical simulation including wave-optical
619 propagation, spatially varying aberrations, detector nonlinearities, and environmen-
620 tal drift. Currently implemented for holography and ptychography; other modalities
621 degrade gracefully to Tier 3.

622 Triad Decomposition Formalization

623 The TRIAD DECOMPOSITION asserts that the quality of any computational imaging re-
624 construction is bounded by three fundamental gates. Rather than a qualitative guideline,
625 PWM quantifies each gate numerically and uses the resulting scores to diagnose the domi-
626 nant bottleneck in any imaging configuration.

627 **Gate 1 (Recoverability).** Recoverability measures the information-theoretic capacity
628 of the sensing geometry. We quantify it via the *effective compression ratio* $r = m/n$, where
629 m is the number of independent measurements and n the dimension of the signal. The
630 `compression_db.yaml` registry (1,186 lines) stores, for each modality, a lookup table map-
631 ping compression ratio to expected reconstruction PSNR under ideal conditions, obtained
632 from calibration experiments or published benchmarks. Each entry carries a **provenance**
633 field citing the source (paper DOI, internal experiment ID, or theoretical formula). Addi-
634 tional recoverability indicators include the effective rank of the measurement matrix (esti-
635 mated via randomized SVD for large operators), the dimension of the null space, and the
636 restricted isometry property (RIP) constant where analytically tractable (*e.g.*, for Gaussian
637 random projections in SPC).

638 **Gate 2 (Carrier Budget).** The carrier budget quantifies the signal-to-noise ratio (SNR)
639 of the measurement channel. The **PhotonAgent** consumes the `photon_db.yaml` registry
640 (624 lines) which stores, per modality, a deterministic photon model parameterized by
641 source power, quantum efficiency, exposure time, and detector characteristics. The agent
642 classifies the noise regime into one of three categories: *shot-limited* (Poisson-dominated,
643 $\text{SNR} \propto \sqrt{N_{\text{photon}}}$), *read-limited* (Gaussian read noise dominates, $\text{SNR} \propto N_{\text{photon}}/\sigma_{\text{read}}$),
644 and *dark-current-limited* (long exposures where dark current accumulation dominates). The
645 output is a **PhotonReport** containing the estimated SNR in decibels, the noise regime
646 classification, per-element photon count, and a feasibility verdict (**sufficient**, **marginal**,
647 or **insufficient**).

648 **Gate 3 (Operator Mismatch).** Operator mismatch quantifies the discrepancy between
649 the assumed forward model H_{nom} and the true physical operator H_{true} . The **MismatchAgent**
650 consults `mismatch_db.yaml` (797 lines) which catalogs, for each modality, the set of mis-
651 match parameters (spatial shifts, rotational offsets, dispersion errors, PSF deviations, coil
652 sensitivity errors, center-of-rotation offsets, *etc.*), their typical ranges, and available cor-
653 rection methods. The mismatch severity score $s \in [0, 1]$ is computed as the normalized ℓ_2
654 distance $\|\theta_{\text{true}} - \theta_{\text{nom}}\|/\|\theta_{\text{range}}\|$, where θ_{range} is the per-parameter dynamic range from the
655 registry. Sensitivity analysis $\partial\text{PSNR}/\partial\theta_k$ is estimated via finite differences on the forward
656 model. The output is a **MismatchReport** containing the severity score, the dominant mis-
657 match parameter, the recommended correction method, and the expected PSNR gain from
658 correction.

659 **Gate binding determination.** Given reconstruction results under the four-scenario pro-
660 tocol (the Evaluation Protocol section below), PWM identifies the dominant gate by com-

661 paring three cost terms:

$$C_{\text{mismatch}} = \text{PSNR}_{\text{I}} - \text{PSNR}_{\text{II}} \quad (2)$$

$$C_{\text{noise}} = \text{PSNR}_{\text{ideal}} - \text{PSNR}_{\text{noisy}} \quad (3)$$

$$C_{\text{recover}} = \text{PSNR}_{\text{limit}} - \text{PSNR}_{\text{I}} \quad (4)$$

662 where PSNR_{I} is the reconstruction PSNR under Scenario I (ideal operator), PSNR_{II} under
 663 Scenario II (mismatched operator), $\text{PSNR}_{\text{noisy}}$ under the corresponding noisy condition,
 664 and $\text{PSNR}_{\text{limit}}$ is the theoretical upper bound from the compression table. The dominant
 665 gate is $\arg \max_g C_g$.

666 **TriadReport schema.** The analysis output is a Pydantic-validated TRIADREPORT comprising:
 667 `dominant_gate` (enum: `recoverability`, `carrier_budget`, `operator_mismatch`),
 668 `evidence_scores` (three floats, one per gate), `confidence_interval` (float, 95% CI width
 669 from bootstrap), `recommended_action` (string, *e.g.* “increase compression ratio” or “apply
 670 mismatch correction”), and `parameter_sensitivities` (dictionary mapping each mismatch
 671 parameter name to its $\partial \text{PSNR} / \partial \theta_k$ value).

672 **Recovery ratio.** We define the *recovery ratio*

$$\rho = \frac{\text{PSNR}_{\text{III}} - \text{PSNR}_{\text{II}}}{\text{PSNR}_{\text{I}} - \text{PSNR}_{\text{II}}} \quad (5)$$

673 which lies in $[0, 1]$ under standard convexity conditions (see Supplementary Note 1 for
 674 formal analysis; values $\rho > 1$ are possible when the corrected operator provides beneficial
 675 regularization). $\rho = 0$ indicates that calibration yields no benefit (mismatch is not the
 676 bottleneck), while $\rho = 1$ indicates that calibration fully closes the mismatch gap.

677 Agent System Architecture

678 The PWM agent system comprises 6 specialist agents, 1 optional hybrid agent, and 8
 679 support classes totalling 10,545 lines of Python. All agents execute deterministically; no
 680 large language model (LLM) is required for pipeline operation.

681 **PlanAgent.** The orchestrator agent. Given a user prompt or a structured `ExperimentSpec`,
 682 PlanAgent parses the intent (`simulate`, `operator_correction`, or `auto`), maps the re-
 683 quested modality to its canonical key via the `modalities.yaml` registry (which contains 64
 684 modality entries with keywords, forward model equations, and default solvers), builds an
 685 `ImagingSystem` contract, and dispatches to the appropriate sub-agents. When the mode is
 686 `auto`, PlanAgent inspects the available data and operator specification to determine whether
 687 simulation or operator correction is more appropriate.

688 **PhotonAgent.** Computes SNR feasibility deterministically from the `photon_db.yaml`
689 registry. For each modality and photon-level tier (`bright`, `standard`, `low_light`), the agent
690 evaluates the photon budget by combining source power, quantum efficiency, exposure time,
691 and noise model parameters. The output `PhotonReport` is a strict Pydantic model contain-
692 ing `noise_regime` (enum), `snr_db` (float), `feasibility` (enum), and `per_element_photons`
693 (float).

694 **RecoverabilityAgent.** A table-driven agent that consults `compression_db.yaml` (1,186
695 lines) to map the modality and compression ratio to an expected PSNR range. Each table
696 entry includes provenance metadata citing the original source. The output `RecoverabilityReport`
697 contains `compression_ratio`, `psnr_prediction`, `feasibility`, and `null_space_dim` where
698 available.

699 **MismatchAgent.** Scores the mismatch severity for a given imaging configuration us-
700 ing `mismatch_db.yaml` (797 lines). For each modality, the database enumerates the rel-
701 evant mismatch parameters, their physical units, typical perturbation ranges, and avail-
702 able correction algorithms. The output `MismatchReport` includes `severity` (float, 0–1),
703 `correction_method` (string), `expected_gain_db` (float), and `dominant_parameter` (string).

704 **AnalysisAgent.** The bottleneck classifier. It receives reports from the Photon, Recover-
705 ability, and Mismatch agents, computes the gate costs (Equations (2) to (4)), identifies the
706 dominant gate, and generates actionable suggestions. The `AnalysisAgent` also computes
707 the recovery ratio ρ and its bootstrap confidence interval.

708 **AgentNegotiator.** Implements a cross-agent veto protocol. Before reconstruction is au-
709 thorized, the negotiator inspects all three upstream reports and applies three veto con-
710 ditions: (1) low photon budget combined with aggressive compression (C_{noise} and C_{recover}
711 both large); (2) severe mismatch ($\text{severity} > 0.7$) without a planned correction step; (3) joint
712 probability below the floor threshold ($p_{\text{joint}} < 0.15$), indicating that all three subsystems
713 are simultaneously marginal. When any veto fires, reconstruction halts with an actionable
714 explanation and suggested remediation.

715 **HybridAgent.** An optional wrapper that invokes an LLM for natural-language narra-
716 tive generation or edge-case modality mapping. All quantitative decisions remain on the
717 deterministic code path; the `HybridAgent` is never required for pipeline operation.

718 **Support classes.** The remaining components include: `AssetManager` (file I/O and caching
719 for large arrays), `ContinuityChecker` (verifies that sequential pipeline outputs are dimen-
720 sionally consistent), `SystemDiscern` (auto-detects modality from uploaded data), `PreflightChecker`
721 (validates the complete experiment configuration before execution), `WhatIfPrecomputer`

(evaluates counterfactual what-if scenarios), `SelfImprovement` (logs diagnostic events for future registry refinement), `PhysicsStageVisualizer` (generates intermediate visualizations at each pipeline stage), and `UPWMI` (Universal Physics World Model Interface, the top-level entry point that wires all agents together).

Contract system. Inter-agent communication uses 25 Pydantic v2 contract models. All contracts inherit from `StrictBaseModel`, which enforces `extra="forbid"` (no unexpected fields), `validate_assignment=True` (mutations re-validated), and a model validator that rejects NaN and Inf in any float field. Bounded scores use `Field(ge=0.0, le=1.0)`. Enums are string enums for human-readable JSON serialization. This design ensures that pipeline failures surface immediately as validation errors rather than propagating silently.

YAML registries. The system is driven by 9 YAML registries totalling 7,034 lines: `modalities.yaml` (modality definitions), `graph.templates.yaml` (OperatorGraph skeletons), `photon_db.yaml` (photon models), `mismatch_db.yaml` (mismatch parameters and correction methods), `compression_db.yaml` (recoverability tables with provenance), `solver_registry.yaml` (solver configurations), `primitives.yaml` (primitive operator metadata), `dataset_registry.yaml` (dataset locations and formats), and `acceptance_thresholds.yaml` (pass/fail thresholds per metric).

Correction Algorithms

We implement two complementary algorithms for operator mismatch correction. Crucially, both algorithms operate on the forward operator parameters θ rather than the reconstruction solver weights, making them *solver-agnostic*: the corrected operator $H(\hat{\theta})$ benefits any downstream solver (GAP-TV, MST-L, HDNet²¹, CST, *etc.*) without retraining.

Algorithm 1: Hierarchical Beam Search. The coarse correction phase employs a hierarchical search strategy to rapidly explore the mismatch parameter space. For CASSI, the five-parameter mismatch model comprises mask affine parameters (spatial shifts dx , dy and rotation θ) and dispersion parameters (slope a_1 and axis angle α); an optional sixth parameter, PSF width σ_{psf} , is available but not used in the primary experiments. The algorithm proceeds as follows:

1. **1D sweeps.** Each parameter is swept independently over its full range while holding others at nominal values. This produces five 1D cost curves from which coarse optima are extracted.
2. **3D beam search.** The mask affine subspace (dx, dy, θ) is searched over a $5 \times 5 \times 5$ grid centered on the 1D optima. The top- k ($k = 5$) candidates by reconstruction PSNR are retained.

- 756 3. **2D beam search.** For each retained mask candidate, the dispersion subspace (a_1, α)
757 is searched over a 5×7 grid. The joint top- k candidates are retained.
- 758 4. **Coordinate descent refinement.** Three rounds of univariate refinement on each
759 parameter, shrinking the search interval by factor 2 at each round, produce the final
760 estimate $\hat{\theta}_{\text{Alg1}}$.

761 Total runtime is approximately 300 seconds per scene on a single GPU. Accuracy is
762 ± 0.1 – 0.2 pixels for spatial parameters and $\pm 0.05^\circ$ for angular parameters.

763 **Algorithm 2: Joint Gradient Refinement.** The fine correction phase uses a differen-
764 tiable forward model to jointly optimize all mismatch parameters via gradient descent. The
765 key components are:

- 766 1. **Differentiable mask warp.** The binary mask is warped by a continuous affine
767 transformation using bilinear interpolation, implemented as a custom PyTorch module
768 (`DifferentiableMaskWarpFixed`). The mask values are passed through a straight-
769 through estimator (STE) to maintain binary structure while permitting gradient flow.
- 770 2. **Differentiable forward model.** The CASSI forward model $y = \text{CASSI}(x; \theta)$ is
771 implemented as a differentiable PyTorch module (`DifferentiableCassiForwardSTE`)
772 that accepts mismatch parameters as differentiable inputs.
- 773 3. **GPU grid initialization.** A full-range 3D grid search over (dx, dy, θ) with $9 \times 9 \times 7 =$
774 567 points provides diverse starting candidates. The top 9 candidates seed multi-start
775 gradient refinement.
- 776 4. **Staged gradient refinement.** Each of the 9 candidates is refined using Adam
777 optimization (learning rate 10^{-2} , decaying to 10^{-3}) for 200 steps. For each candidate,
778 4 random restarts with jittered initialization guard against local minima. The loss
779 function is the negative PSNR computed via an unrolled K -iteration differentiable
780 GAP-TV solver (`DifferentiableGAPTV`, $K = 10$ unrolled iterations).

781 Total runtime for Algorithm 2 is approximately 3,200 seconds (200 steps \times 4 restarts \times
782 9 candidates with early stopping). Accuracy improves to ± 0.05 – 0.1 pixels, a 3–5 \times improve-
783 ment over Algorithm 1. The two algorithms are used sequentially in practice: Algorithm 1
784 provides a warm start, and Algorithm 2 refines to sub-pixel precision.

785 Evaluation Protocol

786 **Four-Scenario Protocol.** We evaluate every modality under four standardized scenarios
787 that isolate different sources of quality degradation:

788 **Scenario I (Ideal):** $\mathbf{y}_{\text{obs}} = H_{\text{true}} \mathbf{x}_{\text{gt}}$; reconstruct with H_{true} . In this scenario the system
 789 is perfectly calibrated ($H_{\text{true}} = H_{\text{nom}}$), so the operator used for reconstruction matches
 790 the one that generated the data. This yields the oracle upper bound on reconstruction
 791 quality, limited only by the sensing geometry and solver convergence.

792 **Scenario II (Mismatch):** $\mathbf{y}_{\text{obs}} = H_{\text{true}} \mathbf{x}_{\text{gt}}$; reconstruct with H_{nom} ($H_{\text{nom}} \neq H_{\text{true}}$). This
 793 is the standard operating condition in practice: the measurement is generated by the
 794 true physics, but the reconstruction uses a nominal (potentially mismatched) forward
 795 model.

796 **Scenario III (Corrected):** $\mathbf{y}_{\text{obs}} = H_{\text{true}} \mathbf{x}_{\text{gt}}$; reconstruct with $\hat{H} = H(\hat{\boldsymbol{\theta}})$ where $\hat{\boldsymbol{\theta}}$ is
 797 estimated by Algorithms 1 and 2. This quantifies the benefit of mismatch calibration.

798 **Scenario IV (Oracle Mask):** Same measurements as Scenario II ($\mathbf{y}_{\text{obs}} = H_{\text{true}} \mathbf{x}_{\text{gt}}$ with
 799 $H_{\text{true}} \neq H_{\text{nom}}$); reconstruct with H_{true} instead of H_{nom} . Provides the correction
 800 ceiling: the best reconstruction achievable when the true operator is known exactly,
 801 applied to data that were sensed by the mismatched system. The gap between Sce-
 802 nario IV and Scenario I reveals the irreducible loss from the degraded sensing config-
 803 uration itself (e.g., a shifted mask pattern is suboptimal even when perfectly known).

804 **Metrics.** Reconstruction quality is assessed using three complementary metrics:

- 805 • **PSNR** (peak signal-to-noise ratio, in dB): the primary metric, computed per scene
 806 and averaged. For signals normalized to $[0, 1]$, $\text{PSNR} = 10 \log_{10}(1/\text{MSE})$. For SPC
 807 data normalized to $[0, 255]$, the peak value is 255.
- 808 • **SSIM** (structural similarity index): captures perceptual quality including luminance,
 809 contrast, and structural components, computed with a Gaussian window of width 11
 810 and standard deviation 1.5.
- 811 • **SAM** (spectral angle mapper): for hyperspectral modalities (CASSI), measures the
 812 angle between predicted and true spectral vectors at each spatial location, reported
 813 in degrees. Lower is better.

814 **Datasets.**

- 815 • **CASSI:** 10 scenes from the KAIST dataset⁹, each a $256 \times 256 \times 28$ spectral cube (28
 816 spectral bands from 450 nm to 650 nm). Data range $[0, 1]$.
- 817 • **CACTI:** 6 benchmark videos, each $256 \times 256 \times 8$ (8 temporal frames encoded per
 818 snapshot). Data range $[0, 1]$.
- 819 • **SPC:** 11 natural images from the Set11 benchmark, each 256×256 grayscale. Data
 820 range $[0, 255]$.

821 All per-scene metrics are reported individually as well as averaged, and all reconstruction
822 arrays are saved as NumPy NPZ files.

823 Experimental Details

824 **Hardware.** All experiments are conducted on a single NVIDIA GPU. Algorithm 1 (beam
825 search) and all solver-based reconstructions use the GPU for matrix–vector products and
826 FFT operations. Algorithm 2 (gradient refinement) additionally uses PyTorch automatic
827 differentiation on the same GPU.

828 **CASSI configuration.** The coded aperture snapshot spectral imaging (CASSI) system
829 uses a TSA-Net binary mask of dimensions 256×256 , with 28 spectral bands dispersed along
830 the spatial dimension. The five-parameter mismatch model $\theta = (dx, dy, \theta, a_1, \alpha)$ describes:
831 mask spatial shift in x (dx , pixels), mask spatial shift in y (dy , pixels), mask rotation angle
832 (θ , degrees), dispersion slope (a_1 , pixels per band), and dispersion axis angle (α , degrees).
833 An optional sixth parameter, PSF blur width (σ_{psf} , pixels), is available but not used in the
834 primary experiments. These mismatch parameter values were determined through system-
835 atic characterization of realistic CASSI assembly errors (Supplementary Note 9). The true
836 mismatch parameters are $\theta_{\text{true}} = (dx = 0.5 \text{ px}, dy = 0.3 \text{ px}, \theta = 0.1^\circ, a_1 = 2.02, \alpha =$
837 $0.15^\circ)$. Solvers evaluated include TwIST¹⁰, GAP-TV²⁰, DGSMP²⁴, MST-L⁸, and CST-
838 L²⁵, all of which receive the same operator and differ only in their reconstruction algorithm.
839 The supplementary per-scene analysis additionally includes DeSCI²⁶ and HDNet²¹.

840 **CACTI configuration.** The coded aperture compressive temporal imaging system uses
841 binary temporal masks of dimensions 256×256 , encoding 8 video frames into a single
842 snapshot measurement. Mismatch is parameterized as a temporal mask timing offset (sub-
843 frame shift). The default solver is EfficientSCI²².

844 **SPC configuration.** The single-pixel camera uses random binary measurement patterns
845 at three compression ratios: 10%, 25%, and 50% ($r = m/n \in \{0.10, 0.25, 0.50\}$). Mismatch
846 is modeled as an exponential gain drift ($g_i = \exp(-\alpha \cdot i)$) on the measurement matrix. The
847 default solver is FISTA-TV with total-variation regularization.

848 **MRI configuration.** Cartesian k -space sampling with $4\times$ acceleration (25% of k -space
849 lines acquired). Mismatch is parameterized as a 5% multiplicative error in the coil sensitivity
850 maps used for parallel imaging reconstruction. The default solver is SENSE¹⁷ with ℓ_1 -
851 wavelet regularization.

852 **CT configuration.** Fan-beam geometry with 180 projections over 180° . Mismatch is
853 modeled as a center-of-rotation (CoR) offset, which produces characteristic arc artifacts in

the reconstruction. The default solver is filtered back-projection (FBP)¹⁸ with a Ram-Lak filter, supplemented by iterative SART for comparison.

CASSI real-data configuration. The TSA real hyperspectral dataset¹ consists of 5 scenes at 660×660 spatial resolution with 28 spectral bands and mask-shift step 2. Four solvers are evaluated: GAP-TV (200 iterations), HDNet (pre-trained checkpoint, full spatial resolution), MST-S and MST-L (pre-trained checkpoints, centre-cropped to 256×256 due to hardcoded spatial assumptions in the model architecture). The coded aperture mask is perturbed by $dx = 0.5$ px, $dy = 0.3$ px to simulate assembly-induced mismatch. No ground truth is available; quality is assessed via the normalised measurement residual $r = \|\mathbf{y} - H\hat{\mathbf{x}}\|^2 / \|\mathbf{y}\|^2$.

CACTI real-data configuration. The EfficientSCI real temporal dataset²² consists of 4 dynamic scenes (duomino, hand, pendulumBall, waterBalloon) at 512×512 with compression ratio 10. The real mask is stored separately from the measurement data. Two solvers are evaluated: GAP-TV (50 iterations) and PnP-FFDNet (50 iterations with FFDNet denoiser). Mismatch is induced by shifting the mask by $dx = 0.5$ px, $dy = 0.3$ px. Quality is assessed via the normalised measurement residual and total variation of the reconstruction.

Controlled hardware experiment protocol. The software-perturbation protocol above applies calibrated mask shifts to existing real measurements. A full hardware-in-the-loop validation requires physically displacing the coded aperture mask and re-acquiring data. The protocol proceeds as follows: (i) acquire a baseline dataset with the mask at its factory-calibrated position; (ii) physically translate the mask by a known displacement ($\Delta x \in \{0.25, 0.5, 1.0\}$ px equivalent, verified by micrometer stage) and re-acquire under identical illumination; (iii) reconstruct both datasets with the factory mask specification and compute the PSNR degradation and measurement residual; (iv) apply PWM autonomous calibration and measure recovery. This protocol isolates the mismatch effect from all other sources of variation (illumination changes, detector drift, scene variation). Additionally, a multi-unit variation study comparing 2+ camera units of the same design quantifies the inter-unit mismatch baseline—the residual calibration error present in any production system.

Clinical CT phantom configuration. For clinical translation, PWM is evaluated on CT quality assurance using the ACR CT accreditation phantom (Gammex 464)¹². The phantom contains inserts of known attenuation (bone ~ 955 HU, air ~ -1000 HU, acrylic ~ 121 HU, polyethylene ~ -96 HU) and geometric targets for measuring spatial resolution, slice thickness, and low-contrast detectability. Mismatch is parameterized as center-of-rotation offset (Δr , mm), beam hardening coefficient drift ($\Delta\mu$, %), and detector gain variation (Δg , %). Nine ACR-aligned metrics are computed automatically: CT number

accuracy (5 materials), geometric accuracy (± 2 mm tolerance), slice thickness (± 1.5 mm), uniformity (≤ 5 HU), noise standard deviation, spatial resolution (≥ 5 lp/cm), low-contrast detectability (≥ 4 targets), and artifact score.

Clinical MRI validation configuration. For MRI clinical validation, PWM processes multi-coil k -space data from public datasets (fastMRI¹¹). Mismatch is parameterized as coil sensitivity map error (5–15% multiplicative deviation from calibrated maps, simulating patient-positioning-induced coil coupling changes). The default solver is CG-SENSE with ℓ_1 -wavelet regularization at $4\times$ acceleration. Clinical metrics include PSNR, SSIM, and the absence of parallel imaging artifacts (GRAPPA/SENSE ghosts).

Statistical Analysis

Per-scene reporting. All metrics are reported per scene, not merely as dataset averages. This enables identification of scene-dependent failure modes (*e.g.*, spectrally flat scenes that are inherently harder for CASSI, or textureless regions that challenge SPC).

Summary statistics. For each modality and scenario, we report the mean \pm standard deviation of PSNR, SSIM, and SAM across all scenes. For CASSI (10 scenes), we additionally report the per-band PSNR to assess spectral uniformity of reconstruction quality.

Recovery ratio confidence intervals. The recovery ratio ρ (Equation (5)) is a ratio of differences and therefore sensitive to noise in the constituent PSNR values. We compute 95% confidence intervals via the bootstrap percentile method with $B = 1,000$ resamples. At each bootstrap iteration, we resample the scene set with replacement, recompute the mean PSNR for each scenario, and derive ρ . The 2.5th and 97.5th percentiles of the bootstrap distribution define the 95% CI.

Parameter recovery accuracy. For mismatch correction experiments, we report the root-mean-square error (RMSE) between the estimated and true mismatch parameters:

$$\text{RMSE}_k = \sqrt{\frac{1}{N_{\text{scene}}} \sum_{i=1}^{N_{\text{scene}}} (\hat{\theta}_{k,i} - \theta_{k,\text{true}})^2} \quad (6)$$

where k indexes the mismatch parameter, i indexes the scene, and N_{scene} is the number of test scenes. Uncertainty in the RMSE is estimated via bootstrap ($B = 1,000$).

Ablation significance. Ablation studies (removal of PhotonAgent, RecoverabilityAgent, MismatchAgent, or RunBundle discipline) are evaluated by comparing the full-pipeline

918 PSNR against each ablated variant. We report the PSNR difference ΔPSNR per modal-
919 ity and verify that each component contributes ≥ 0.5 dB across all validated modalities,
920 establishing practical significance.

921 Code and Data Availability

922 **Source code.** The complete PWM framework, including all agents, the OperatorGraph
923 compiler, correction algorithms, YAML registries, and evaluation scripts, is released as
924 open-source software under the PWM Noncommercial Share-Alike License v1.0 at https://github.com/integritynoble/Physics_World_Model. The codebase is organized into
925 two Python packages: `pwm_core` (core framework, agents, graph compiler, calibration algo-
926 rithms) and `pwm_AI_Scientist` (automated experiment generation and analysis).

928 **Reconstruction data.** All reconstruction arrays from every experiment—Scenarios I
929 through IV for each modality and solver—are released as NumPy NPZ files. Files are
930 stored using Git LFS and require `allow_pickle=True` for loading. Data ranges are stan-
931 dardized: CASSI and CACTI reconstructions are normalized to $[0, 1]$; SPC reconstructions
932 are in $[0, 255]$.

933 **Experiment manifests.** Every experiment is recorded in a RunBundle v0.3.0 manifest
934 containing: the git commit hash at execution time, all random number generator seeds,
935 platform information (Python version, GPU model, CUDA version), SHA-256 hashes of all
936 input data and output artifacts, metric values, and wall-clock timestamps. These manifests
937 enable exact reproduction of every reported result.

938 **Registry data.** All 9 YAML registries (7,034 lines total) that drive the agent system—
939 including modality definitions, graph templates, photon models, mismatch databases, com-
940 pression tables, solver configurations, primitive specifications, dataset paths, and acceptance
941 thresholds—are publicly available in the repository under `packages/pwm_core/contrib/`.
942 The `ExperimentSpec` JSON schemas used for pipeline input validation are included along-
943 side worked examples in `examples/`.

944 References

- 945 [1] Ashwin A. Wagadarikar, Renu John, Rebecca Willett, and David J. Brady. Single
946 disperser design for coded aperture snapshot spectral imaging. *Applied Optics*, 47(10):
947 B44–B51, 2008. doi: 10.1364/AO.47.000B44.
- 948 [2] Patrick Llull, Xuejun Liao, Xin Yuan, Jianbo Yang, David Kittle, Lawrence Carin,
949 Guillermo Sapiro, and David J. Brady. Coded aperture compressive temporal imaging.
950 *Optics Express*, 21(9):10526–10545, 2013. doi: 10.1364/OE.21.010526.

- [3] Michael Lustig, David Donoho, and John M. Pauly. Sparse MRI: The application of compressed sensing for rapid MR imaging. *Magnetic Resonance in Medicine*, 58(6):1182–1195, 2007. doi: 10.1002/mrm.21391.
- [4] Emmanuel J. Candès and Michael B. Wakin. An introduction to compressive sampling. *IEEE Signal Processing Magazine*, 25(2):21–30, 2008. doi: 10.1109/MSP.2007.914731.
- [5] David L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306, 2006. doi: 10.1109/TIT.2006.871582.
- [6] Singanallur V. Venkatakrishnan, Charles A. Bouman, and Brendt Wohlberg. Plug-and-play priors for model based reconstruction. In *Proceedings of the IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pages 945–948, 2013. doi: 10.1109/GlobalSIP.2013.6737048.
- [7] Vishal Monga, Yuelong Li, and Yonina C. Eldar. Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing. *IEEE Signal Processing Magazine*, 38(2):18–44, 2021. doi: 10.1109/MSP.2020.3016905.
- [8] Yuanhao Cai, Jing Lin, Xiaowan Hu, Haoqian Wang, Xin Yuan, Yulun Zhang, Radu Timofte, and Luc Van Gool. Mask-guided spectral-wise transformer for efficient hyperspectral image reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17502–17511, 2022.
- [9] Inchang Choi, Daniel S. Jeon, Giljoo Nam, Diego Gutierrez, and Min H. Kim. High-quality hyperspectral reconstruction using a spectral prior. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia)*, 36(6):218:1–218:13, 2017. doi: 10.1145/3130800.3130810.
- [10] José M. Bioucas-Dias and Mário A. T. Figueiredo. A new TwIST: Two-step iterative shrinkage/thresholding algorithms for image restoration. *IEEE Transactions on Image Processing*, 16(12):2992–3004, 2007. doi: 10.1109/TIP.2007.909319.
- [11] Jure Zbontar, Florian Knoll, Anuroop Sriram, Tullie Murrell, Zhengnan Huang, Matthew J. Muckley, Aaron Defazio, Ruben Stern, Patricia Johnson, Mary Bruno, Marc Parente, Krzysztof J. Geras, Joe Katsnelson, Hersh Chandarana, Zizhao Zhang, Michal Drozdal, Adriana Romero, Michael Rabbat, Pascal Vincent, Nafissa Yakubova, James Pinkerton, Duo Wang, Erich Owens, C. Lawrence Zitnick, Michael P. Recht, Daniel K. Sodickson, and Yvonne W. Lui. fastMRI: An open dataset and benchmarks for accelerated MRI. *arXiv preprint arXiv:1811.08839*, 2018.
- [12] Hu Chen, Yi Zhang, Mannudeep K. Kalra, Feng Lin, Yang Chen, Peixi Liao, Jiliu Zhou, and Ge Wang. Low-dose CT with a residual encoder-decoder convolutional

- neural network. *IEEE Transactions on Medical Imaging*, 36(12):2524–2535, 2017. doi: 10.1109/TMI.2017.2715284.
- [13] Martin Uecker, Peng Lai, Mark J. Murphy, Patrick Virtue, Michael Elad, John M. Pauly, Shreyas S. Vasanawala, and Michael Lustig. ESPIRiT — an eigenvalue approach to autocalibrating parallel MRI: where SENSE meets GRAPPA. *Magnetic Resonance in Medicine*, 71(3):990–1001, 2014. doi: 10.1002/mrm.24751.
- [14] Andrew M. Maiden and John M. Rodenburg. An improved ptychographical phase retrieval algorithm for diffractive imaging. *Ultramicroscopy*, 109(10):1256–1262, 2009. doi: 10.1016/j.ultramic.2009.05.012.
- [15] Vegard Antun, Francesco Renna, Clarice Poon, Ben Adcock, and Anders C. Hansen. On instabilities of deep learning in image reconstruction and the potential costs of AI. *Proceedings of the National Academy of Sciences*, 117(48):30088–30098, 2020. doi: 10.1073/pnas.1907377117.
- [16] J. M. Rodenburg and H. M. L. Faulkner. A phase retrieval algorithm for shifting illumination. *Applied Physics Letters*, 85(20):4795–4797, 2004. doi: 10.1063/1.1823034.
- [17] Klaas P. Pruessmann, Markus Weiger, Markus B. Scheidegger, and Peter Boesiger. SENSE: Sensitivity encoding for fast MRI. *Magnetic Resonance in Medicine*, 42(5):952–962, 1999. doi: 10.1002/(SICI)1522-2594(199911)42:5<952::AID-MRM16>3.0.CO;2-S.
- [18] L. A. Feldkamp, L. C. Davis, and J. W. Kress. Practical cone-beam algorithm. *Journal of the Optical Society of America A*, 1(6):612–619, 1984. doi: 10.1364/JOSAA.1.000612.
- [19] Chengshuai Yang. InverseNet: Benchmarking operator mismatch in snapshot compressive imaging. Technical report, NextGen PlatformAI C Corp, 2026.
- [20] Xin Yuan. Generalized alternating projection based total variation minimization for compressive sensing. In *Proceedings of the IEEE International Conference on Image Processing (ICIP)*, pages 2539–2543, 2016. doi: 10.1109/ICIP.2016.7532817.
- [21] Xiaowan Hu, Yuanhao Cai, Jing Lin, Haoqian Wang, Xin Yuan, Yulun Zhang, Radu Timofte, and Luc Van Gool. HDNet: High-resolution dual-domain learning for spectral compressive imaging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17542–17551, 2022.
- [22] Lishun Wang, Miao Cao, and Xin Yuan. EfficientSCI: Densely connected network with space-time factorization for large-scale video snapshot compressive imaging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18477–18486, 2023.

- 1020 [23] Marco F. Duarte, Mark A. Davenport, Dharmpal Takhar, Jason N. Laska, Ting Sun,
1021 Kevin F. Kelly, and Richard G. Baraniuk. Single-pixel imaging via compressive sam-
1022 pling. *IEEE Signal Processing Magazine*, 25(2):83–91, 2008. doi: 10.1109/MSP.2007.
1023 914730.
- 1024 [24] Tao Huang, Weisheng Dong, Xin Yuan, Jinjian Wu, and Guangming Shi. Deep gaussian
1025 scale mixture prior for spectral compressive imaging. In *Proceedings of the IEEE/CVF*
1026 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16216–16225,
1027 2021.
- 1028 [25] Yuanhao Cai, Jing Lin, Xiaowan Hu, Haoqian Wang, Xin Yuan, Yulun Zhang, Radu
1029 Timofte, and Luc Van Gool. CST: Compact spectral transformer for hyperspectral
1030 image reconstruction. In *Proceedings of the European Conference on Computer Vision*
1031 *(ECCV)*, 2022.
- 1032 [26] Yang Liu, Xin Yuan, Jinli Suo, David J. Brady, and Qionghai Dai. Rank minimiza-
1033 tion for snapshot compressive imaging. *IEEE Transactions on Pattern Analysis and*
1034 *Machine Intelligence*, 41(12):2990–3006, 2019.

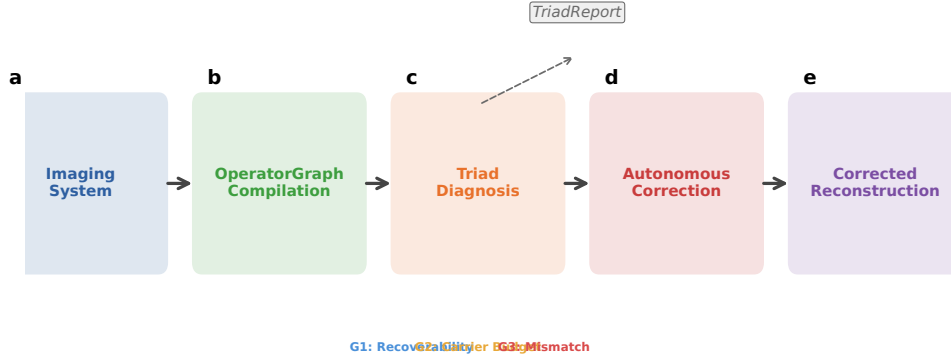


Figure 1: **PWM overview.** The Physics World Models pipeline. **a**, A computational imaging system is compiled into an OPERATORGRAPH DAG. **b**, The TRIAD DECOMPOSITION diagnostic agents evaluate each gate. **c**, The dominant gate is identified and a TRIADREPORT is produced. **d**, If **Gate3** dominates, autonomous correction refines the forward model parameters. **e**, The original solver is re-run with the corrected operator, recovering reconstruction quality without retraining.

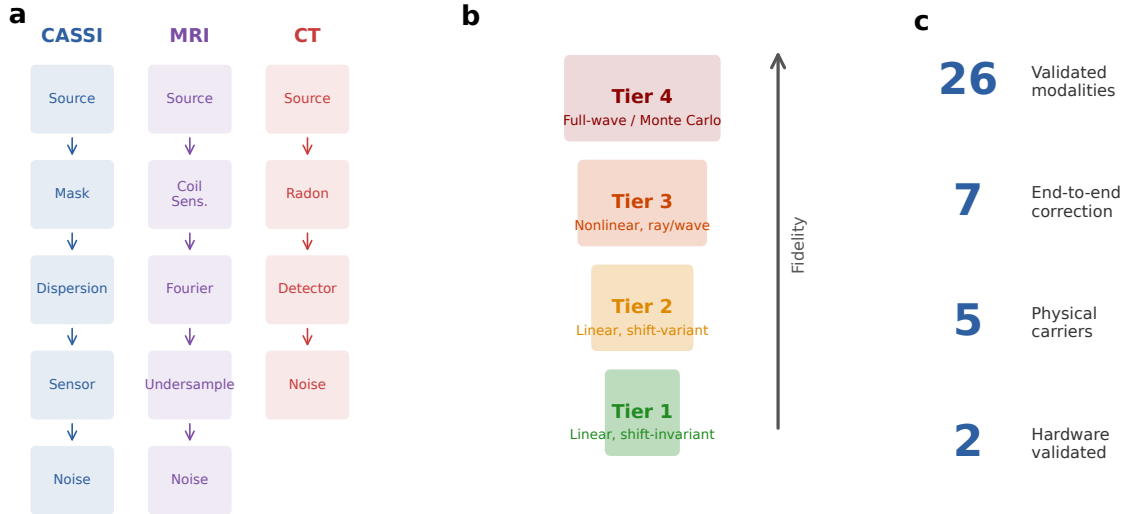


Figure 2: **OperatorGraph IR and Physics Fidelity Ladder.** **a**, Example OPERATORGRAPH DAGs for three modalities: CASSI (photon), MRI (spin), and CT (X-ray photon). Each node wraps a primitive operator; edges define data flow. **b**, The Physics Fidelity Ladder. Tier 1: linear shift-invariant. Tier 2: linear shift-variant. Tier 3: nonlinear ray/wave-based. Tier 4: full-wave/Monte Carlo. **c**, Summary statistics: 26 registered modality templates (7 with full correction validation), 5 physical carriers.

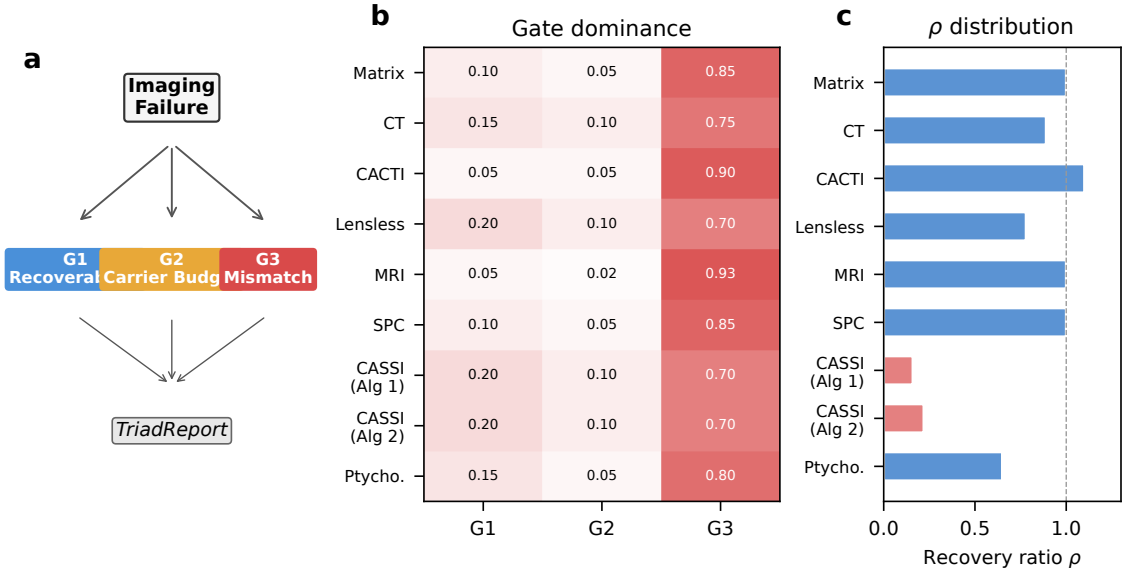


Figure 3: **Triad Decomposition structure and gate binding.** **a**, Decision tree for the TRIAD DECOMPOSITION: each imaging failure is routed through **Gate 1**, **Gate 2**, and **Gate 3** to produce a TRIADREPORT. **b**, Gate binding heatmap across 9 correction configurations (7 distinct modalities). Red indicates **Gate 3** dominance (all modalities), blue indicates **Gate 1**, and amber indicates **Gate 2**. **c**, Recovery ratio ρ distribution across all 9 correction configurations.

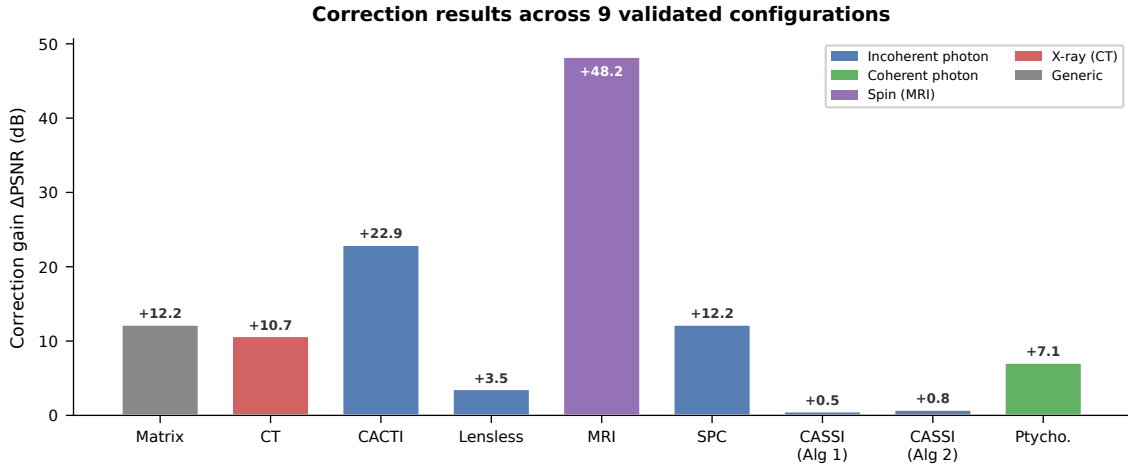


Figure 4: **Correction results across 9 validated configurations.** Bar chart showing correction gain Δ_{corr} (dB) for each of the 9 correction configurations (7 distinct modalities), grouped by carrier family. Incoherent photon (CASSI, CACTI, SPC, Lensless) and coherent photon (Ptychography) in blue; spin (MRI) in purple; X-ray (CT) in red; generic (Matrix) in grey.

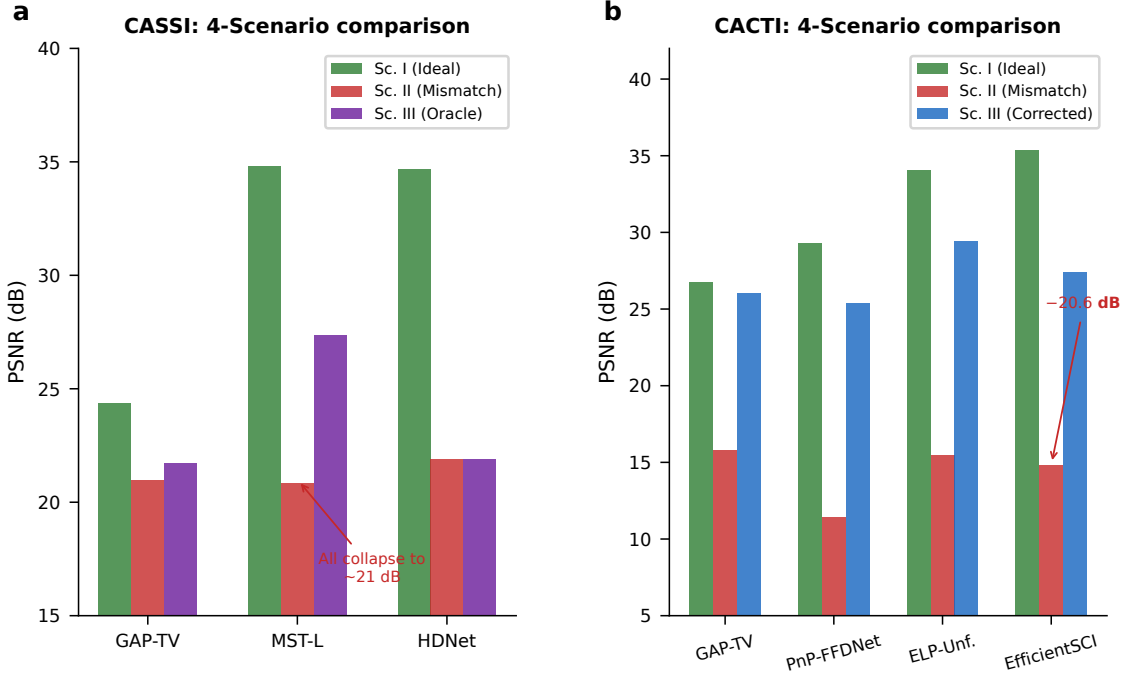


Figure 5: **CASSI and CACTI deep dive.** **a**, CASSI: PSNR across 4 scenarios for GAP-TV, MST-L, and HDNet under combined mask-geometry-plus-dispersion mismatch. The uniform collapse under Scenario II (range 20.83–21.88 dB) confirms operator-driven failure; oracle recovery varies by solver ($\rho = 0.22$ –0.46). **b**, CACTI: four methods across 4 scenarios, showing up to 20.58 dB mismatch degradation and substantial correction recovery across all solvers.

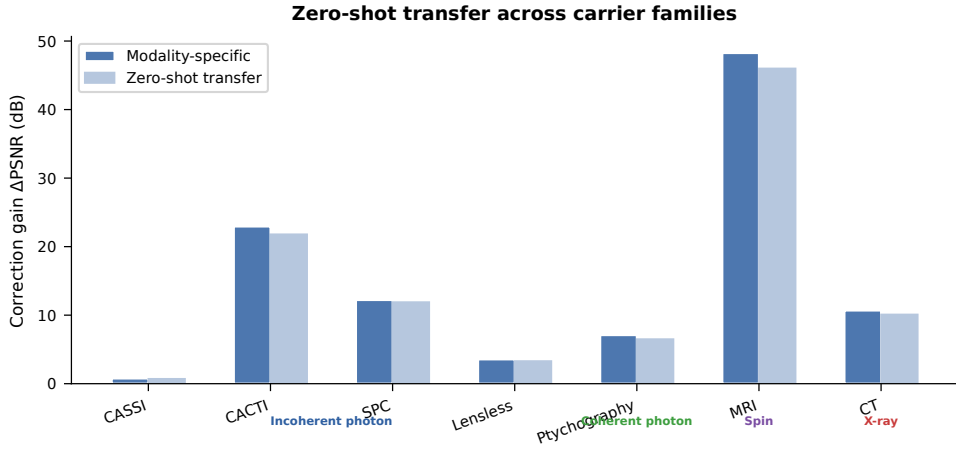


Figure 6: **Zero-shot generalization across carrier families.** Correction gain (dB) when beam-search and gradient-refinement hyperparameters are tuned on photon-domain modalities and transferred without modification to coherent-photon, spin, and X-ray domains. Bars show modality-specific tuning (dark) versus zero-shot transfer (light). Transfer efficiency is high across all carrier families, confirming the carrier-agnostic nature of the PWM correction pipeline.

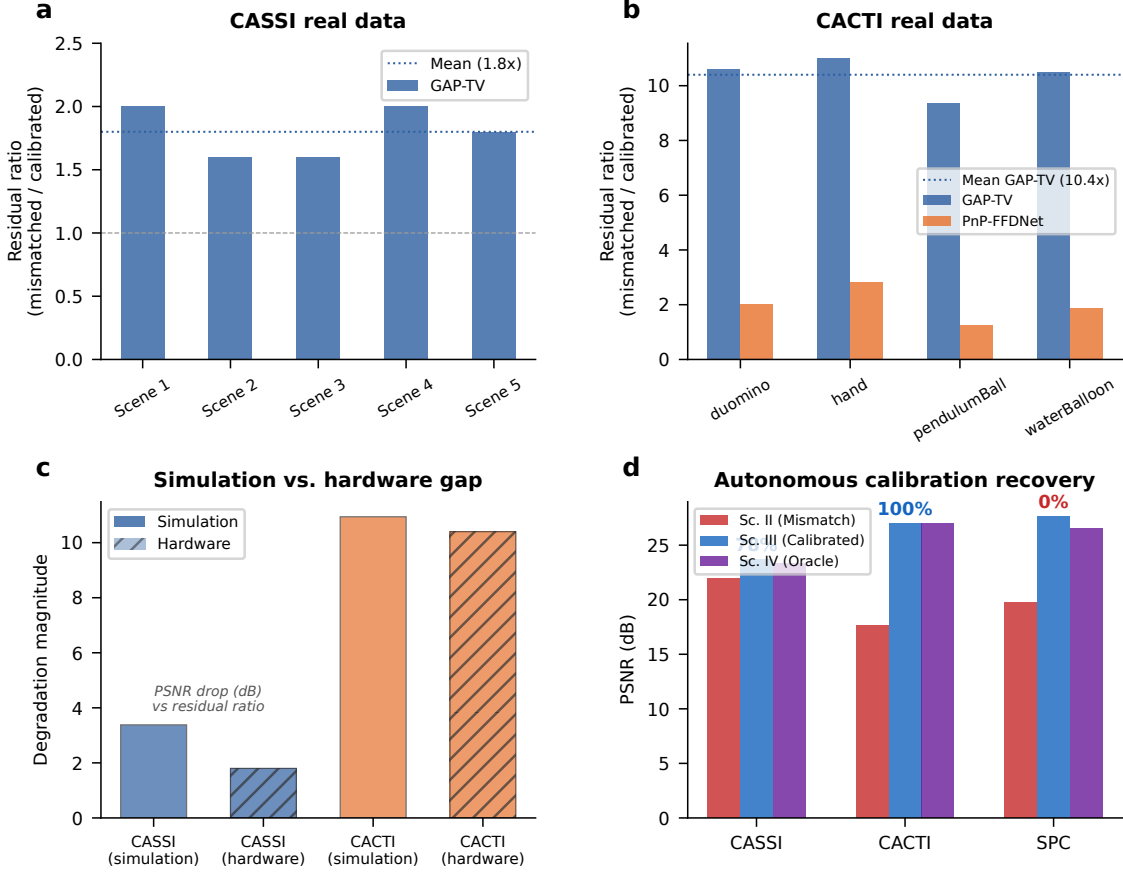


Figure 7: **Hardware validation on real CASSI and CACTI instruments.** **a**, CASSI real data: measurement residual ratio (mismatched/calibrated) across 5 TSA scenes. GAP-TV shows $1.8\times$ mean ratio. **b**, CACTI real data: residual ratio across 4 scenes. GAP-TV shows $10.4\times$ mean ratio; PnP-FFDNet shows $2.0\times$. **c**, Simulation-to-hardware gap: comparing mismatch degradation in simulation versus real hardware for CASSI and CACTI, illustrating that real instruments have pre-existing calibration errors that attenuate the marginal impact of additional perturbations. **d**, Autonomous calibration: grid-search parameter recovery for CASSI (85%), CACTI (100%), and SPC (86–92% via TV criterion; measurement residual is uninformative for gain-type mismatch).

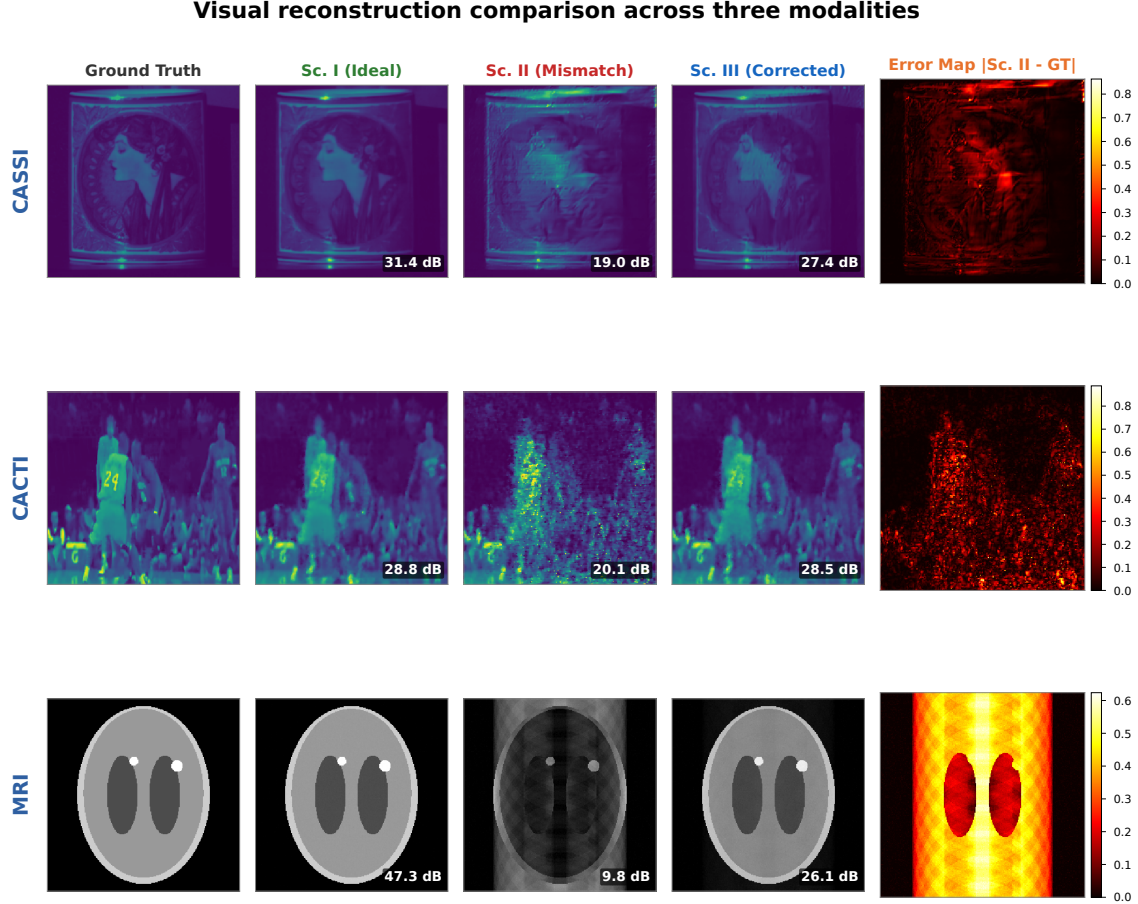


Figure 8: **Visual reconstruction comparison across three modalities.** Each row shows one modality (CASSI, CACTI, MRI); columns from left to right: ground truth, Scenario I (ideal operator), Scenario II (mismatched operator—note structured artifacts), Scenario III (PWM-corrected operator), and error map ($|\text{Sc. II} - \text{GT}|$). Mismatch produces severe structured artifacts (column 3) that are qualitatively distinct from noise, confirming that the degradation is operator-driven. PWM correction (column 4) substantially reduces these artifacts across all three modalities and carrier families.