

The Theory Behind Stellar Mind AI's Application Tools for Space Biology

Quemada-Torres, Eulogio, Román-Sánchez, Alejandro, Montes-Pérez, Javier, Decena-Gómez, Macorís, Dávila-Moreno, Candela

University of Málaga, Blvr. Louis Pasteur, 35, 29071 Málaga, Spain.

To cite this article: Quemada-Torres, E., Román-Sánchez, A., Montes-Pérez, J., Decena-Gómez, M., Dávila-Moreno, C. 2025. The Theory Behind Stellar Mind AI's Application Tools for Space Biology. NASA Space Apps Challenge.

Abstract

We present a formal and fully theoretical specification of *Stellar Mind AI*, a space-biology knowledge engine featuring a graph-centric retrieval-augmented generation (GraphRAG) layer, an *Assay Finder* that maps natural-language intent to a normalized assay manifold, a decision-oriented *Gap Finder* based on interpretable coverage functionals, and a reserved *Mission Builder*. We define objects, invariants, and compact scoring principles grounded in spectral graph theory, submodular optimization, and exponential-family modeling. No implementation or API detail is discussed.

Keywords: Retrieval-Augmented Generation, Knowledge Graphs, Space Biology, Spectral Methods, Submodularity

Notation. Let $G = (V, E)$ be a typed graph; L_H the normalized Laplacian of a subgraph H ; $\lambda_2(L_H)$ its algebraic connectivity; π_q^{PPR} a seeded diffusion from query q ; $\mathbb{I}\{\cdot\}$ the indicator. We write $\langle \cdot, \cdot \rangle$ for inner products, $\text{KL}(\cdot \| \cdot)$ for Kullback–Leibler divergence.

1. Introduction

Space-biology knowledge is fundamentally relational. Transformer encoders model sequence-level semantics (Vaswani et al., 2017), while knowledge graphs canonize entities and typed relations (Hogan et al., 2021). Retrieval-augmented generation (RAG) composes both views (Manning et al., 2008; Gao and et al., 2023), and graph-based RAG (*GraphRAG*) elevates retrieval from documents to structured subgraphs, enabling global sensemaking and multi-hop attribution (Edge et al., 2024; Microsoft Research, 2024; Larson and Truitt, 2024; Peng and et al., 2024; Zhang and et al., 2025). We formalize a theoretical layer for four modules—GraphRAG, Assay Finder, Gap Finder, and a reserved Mission Builder—using spectral connectivity, random-walk diffusion, and submodular selection as first-class primitives (Chung, 1997; Andersen et al., 2006; Tong et al., 2006; Krause and Golovin, 2014).

2. GraphRAG: Subgraph Retrieval and Faithful Conditioning

Objects. From corpus \mathcal{D} we induce a typed, heterogeneous graph $G = (V, E, \tau_V, \tau_E)$ with canonical labels and optional embeddings $\mathbf{z}_v, \mathbf{z}_e$ for textual evidence and mentions. Typed relations encode a constrained hypothesis class whose

latent geometry can be regularized by knowledge-graph embedding priors (e.g., translational, complex, rotational constraints) to promote type-consistent neighborhoods and equivariant scoring (Bordes et al., 2013; Trouillon et al., 2016; Sun et al., 2019). GraphRAG assumes that global sensemaking emerges from *structure-aware retrieval* rather than bag-of-passages: queries target *subgraphs* that preserve multi-hop support and community-level semantics (Edge et al., 2024; Larson and Truitt, 2024; Microsoft Research, 2024).

Intent and constraints. A natural-language query q is mapped to an abstract intent $\mathcal{I}(q) = (\mathcal{C}, \mathcal{M}, \mathcal{F})$: (i) hard constraints \mathcal{C} over entities/types; (ii) soft *motifs* \mathcal{M} (relation templates, e.g., *perturbs* \circ *expresses*); (iii) semantic facets \mathcal{F} (topic/domain vectors). These play the roles of feasibility, weak pattern matching, y *coverage* semántica respectivamente. Given $\mathcal{I}(q)$, retrieval selects a connected subgraph $H \subseteq G$ as the maximizer of the composite objective

$$H^* \in \arg \max_{H \subseteq G} \lambda_1 \Phi_m(H; \mathcal{C}, \mathcal{M}, \mathcal{F}) + \lambda_2 \Psi_{\text{coh}}(H) + \lambda_3 \Xi_{\text{cov}}(H) - \lambda_4 R_{\text{red}}(H), \quad (1)$$

subject to type-compatibility and connectivity. Each term captures a distinct theoretical desideratum:

- *Motif/constraint satisfaction* Φ_m scores typed pattern match with slack on facets, promoting subgraphs whose path-types agree with \mathcal{M} and whose node/edge labels are feasible under \mathcal{C} .
- *Cohesion* $\Psi_{\text{coh}}(H) = \lambda_2(L_H)$ uses the algebraic connectivity (second eigenvalue of the normalized Laplacian L_H). By Cheeger-type inequalities, large λ_2 implies low conductance and internally well-connected

subgraphs (spectral cuts avoid spurious bridges) (Chung, 1997; Shi and Malik, 2000; Ng et al., 2001).

- *Coverage* $\Xi_{\text{cov}}(H) = \sum_{v \in V(H)} \pi_q^{\text{PPR}}(v)$ aggregates a seeded diffusion (personalized PageRank / random-walk-with-restart) from query-induced seeds, thus favoring regions with high stationary mass near relevant anchors (Andersen et al., 2006; Tong et al., 2006).
- *Redundancy penalty* R_{red} discourages near-duplicates and improves diversity, instantiated either as a facility-location repulsion (submodular) or as a DPP energy over features/paths with kernel K on $V(H)$ (Kulesza and Taskar, 2012; Krause and Golovin, 2014).

Why spectral + diffusion? Spectral cohesion (λ_2) controls the isoperimetric profile of H ; diffusion captures *locality* relative to q . Together they select subgraphs that are (i) *internally tight* (few boundary edges), (ii) *externally well-anchored* near query seeds, and (iii) *semantically on-pattern* via Φ_m . This mirrors local partitioning results where PPR “finds” low-conductance sets around a seed (Andersen et al., 2006), while RWR provides relevance scores to rank candidate neighborhoods (Tong et al., 2006).

Algorithmic relaxations and guarantees. The exact problem (1) couples connectivity, pattern constraints y diversidad, lo que es NP-hard. Dos relajaciones cubren casos prácticos:

(R1) *Relajación espectral.* Optimizar Rayleigh quotients $\mathbf{u}^\top L_H \mathbf{u}$ con restricciones suaves de pertenencia produce vectores propios que, al umbralizar, dan cortes de baja conductancia con garantías clásicas (Chung, 1997; Ng et al., 2001; Shi and Malik, 2000). Se puede incorporar Φ_m ponderando L (p.ej., aristas compatibles) o prefiltrando el soporte.

(R2) *Relajación submodular.* Si $\Phi_m + \Xi_{\text{cov}}$ es (aprox.) monótona submodular y R_{red} es un matroide (o se usa un surrogate submodular de DPP), el greedy bajo presupuesto/knapsack obtiene factor $(1 - 1/e)$ (Krause and Golovin, 2014). La conectividad se repara a posteriori vía augmentación de Steiner mínima sobre el corte seleccionado, o se impone *a priori* con restricciones de árbol en expansión.

Faithful conditioning and invariants. Let $\mathcal{P}(H)$ denote typed *minimal support paths*. Conditioning is constrained by three invariants (compact form to fit a column):

- (I) $\forall c \in \mathcal{C}_{\text{at}} \exists p \in \mathcal{P}(H) : p \models c$ (provenance closure)
- (II) $\mathcal{N} \circ \mathcal{N} = \mathcal{N}$ on $V(H), E(H)$ (idempotent normalization)
- (III) $\max_{|\text{ctx}| \leq B} \{\text{DPP}(\text{ctx}) + \text{Cov}(\text{ctx})\}$ (diversity packing)

(I) forces every atomic claim to be backed by a typed path; (II) avoids drifting entity canonicalization across decoding steps; (III) formalizes context selection as a small- B packing that balances diversity (repulsion) and coverage, consistent with DPP and submodular objectives used in diverse summarization (Kulesza and Taskar, 2012; Krause and Golovin, 2014; Yu and et al., 2024). These constraints align with attribution/evaluation practices in RAG and with the “local-to-global” design of GraphRAG (Edge et al., 2024; Larson and Truitt, 2024).

Communities and local-to-global sensemaking. Global queries (e.g., “main themes”) are ill-posed for paragraph retrieval but natural for graph *communities*. Let $\{C_k\}$ be clusters (spectral or modularity-based) (Newman, 2006; Blondel et al., 2008; Ng et al., 2001). GraphRAG precomputes summaries $s_k = \mathcal{S}(C_k)$ and, at query time, weights them by $\lambda_2(L_{C_k})$ (internal coherence) and $\pi_q^{\text{PPR}}(C_k)$ (query affinity), then composes partial answers \rightarrow final synthesis (Edge et al., 2024). This *two-level* pipeline explains empirically observed gains on “global” QFS-style questions while preserving attribution to H ’s paths (Edge et al., 2024; Microsoft Research, 2024).

Relation to standard RAG. Classical RAG scores independent chunks; GraphRAG scores *structured* neighborhoods. Replacing TF-IDF/dot-product retrieval by (1) yields (i) fewer boundary hallucinations (vía λ_2 y caminos tipados), (ii) mejor cobertura multihop (difusión sembrada), y (iii) *diversidad controlada* del contexto (DPP/submodular), lo cual facilita la evaluación de fe y cobertura (Gao and et al., 2023; Yu and et al., 2024).

Design checklist (theoretical). Typed feasibility (respecto a \mathcal{C}); Motif slack calibrado en Φ_m ; Cohesion floor $\lambda_2(L_H) \geq \epsilon$; Coverage budget via PPR mass $\sum_{v \in V(H)} \pi_q^{\text{PPR}}(v) \geq \eta$; Diversity budget $|\text{ctx}| \leq B$ con repulsión DPP o facility-location. Este “perfil” abstrae implementaciones y deja claros los grados de libertad del modelo (Andersen et al., 2006; Tong et al., 2006; Kulesza and Taskar, 2012; Krause and Golovin, 2014).

3. Assay Finder: Exponential-Family Projection on a Normalized Manifold

Manifold and normalization. Let the assay manifold be

$$\mathcal{X} = \mathcal{O} \times \mathcal{T} \times \mathcal{C} \times \mathcal{A},$$

where \mathcal{O} (organism), \mathcal{T} (parent tissue), \mathcal{C} (condition; coarse Spaceflight vs. Ground/Analog), and \mathcal{A} (assay type) are finite controlled sets. A normalization operator \mathcal{N} acts on raw metadata strings and induces an equivalence relation $\sim_{\mathcal{N}}$ whose classes are the canonical tokens. We require

$$\mathcal{N} \circ \mathcal{N} = \mathcal{N} \quad \text{and} \quad x \sim_{\mathcal{N}} x' \Rightarrow \text{same axes in } \mathcal{X},$$

i.e., idempotence and axis-consistent canonicalization. In practice, \mathcal{N} is grounded by anatomy/biomedical ontologies on \mathcal{T} and \mathcal{O} (e.g., parent-tissue lifts and taxon canonicalization) (Mungall and et al., 2012; Smith and et al., 2007). Let $\Omega \subseteq \mathcal{X}$ be the observed (nonempty) subset.

Semantic encoding and exponential family. A request q is mapped to a natural parameter $\theta = \theta(q) \in \mathbb{R}^d$ and to sufficient statistics $\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^d$ (token indicators and semantic features). We define the exponential family

$$\pi_{\theta}(x) = \frac{1}{Z(\theta)} \exp(\langle \theta, \mathbf{f}(x) \rangle), \quad Z(\theta) = \sum_{x \in \mathcal{X}} e^{\langle \theta, \mathbf{f}(x) \rangle},$$

with standard properties of minimal, regular exponential families on finite domains (Brown, 1986; Wainwright and Jordan, 2008). A *concept set* is the superlevel region

$$\mathcal{A}_q(\tau) = \{x \in \mathcal{X} : \pi_{\theta}(x) \geq \tau\}, \quad \tau \in (0, 1).$$

We assume a feasibility polytope $\mathcal{M} \subset \Delta(\mathcal{X})$ enforcing ontology/legal constraints (e.g., type compatibility). A raw intent distribution μ (e.g., a soft prior over \mathcal{X}) is projected as an I -projection

$$\hat{\pi} = \arg \min_{\pi \in \mathcal{E} \cap \mathcal{M}} \text{KL}(\mu \| \pi), \quad \mathcal{E} = \{\pi_\theta : \theta \in \mathbb{R}^d\},$$

which is equivalent to maximum-entropy under moment constraints (Jaynes, 1957; Csiszár, 1975; Amari and Nagaoka, 2000).

Principles. (P1) \mathcal{N} idempotent; (P2) *support discipline*: if selection is restricted to the observed mask (i.e., we enforce $\text{supp}(\pi) \subseteq \Omega$ via \mathcal{M}), then the projection introduces no mass outside Ω ; otherwise, forward KL does not penalize assigning mass where μ has zeros, so additional support may appear (Csiszár, 1975). (P3) *auditability*: membership $x \in \mathcal{A}_q(\tau)$ is justified by typed paths in G (cf. Sec. 2).

Existence/uniqueness and moment matching. Since \mathcal{E} is a regular minimal exponential family on finite \mathcal{X} , the objective is strictly convex in π and the feasible set $\mathcal{E} \cap \mathcal{M}$ is convex; thus $\hat{\pi}$ exists and is unique whenever μ admits a feasible moment vector $\mathbb{E}_\mu[\mathbf{f}]$ in the relative interior of the convex moment set induced by \mathcal{M} (Brown, 1986; Wainwright and Jordan, 2008). Moreover, $\hat{\theta}$ satisfies the moment conditions

$$\mathbb{E}_{\hat{\pi}}[\mathbf{f}] = \mathbb{E}_\mu[\mathbf{f}] \quad (\text{maximum-entropy}/I\text{-projection}).$$

Grouping, deduplication, comparability. We formalize three operations, expressed purely over (\mathcal{X}, Ω) :

(G1) *Grouping by technology.* Let $g : \mathcal{X} \rightarrow \mathcal{A}$ be the projection $g(o, t, c, a) = a$. For any distribution ρ on \mathcal{X} (e.g., $\rho = \hat{\pi}$), define group-mass

$$\Gamma(a) = \sum_{x \in \mathcal{X} : g(x)=a} \rho(x), \quad a \in \mathcal{A}.$$

(G2) *Deduplication.* Let \equiv be an equivalence relation on Ω capturing duplicates at the assay granularity (e.g., same canonical assay-name within dataset). The canonical representative map $D : \Omega \rightarrow \Omega/\equiv$ is idempotent and order-preserving for any scoring that is constant on classes.

(G3) *Comparability flags.* Define binary functionals on $x = \langle o, t, c, a \rangle$:

$$\chi_{\text{Sf}}(x) = \mathbb{1}\{c = \text{Sf}\}, \quad \chi_{\text{Gnd}}(x) = \mathbb{1}\{c = \text{Gnd}\},$$

and a “both” flag at the (o, t, a) level:

$$\chi_{\text{Both}}(o, t, a) = \mathbb{1}\{\exists c \neq c' : \langle o, t, c, a \rangle, \langle o, t, c', a \rangle \in \Omega\}.$$

These flags are measurable in Ω and invariant under \mathcal{N} .

Prioritization on the manifold. Let $\kappa : \mathcal{X} \rightarrow \mathbb{N}$ count instances in Ω (cf. Gap Finder). Define compact, interpretable signals on x :

$$\eta_1(x) = \chi_{\text{Sf}}(x), \quad \eta_2(x) = \chi_{\text{Gnd}}(x), \quad \eta_3(x) = \chi_{\text{Both}}(o, t, a),$$

$$\eta_4(x) = \min\left(1, \frac{\kappa(\langle o, t, \text{Gnd}, a \rangle)}{3}\right)$$

$$\eta_5(x) = \min\left(1, \frac{\kappa(\langle o, t, \text{Sf}, a \rangle)}{3}\right)$$

and a neighborhood density on the product space,

$$\eta_6(x) = \frac{1}{Z_\rho} \sum_{y \in \mathcal{N}_\rho(x)} \kappa(y),$$

with \mathcal{N}_ρ a product-metric neighborhood and Z_ρ a normalizer. For $\alpha \in \mathbb{R}_{\geq 0}^6$, the *assay-score*

$$R(x) = \langle \alpha, \eta(x) \rangle \quad \text{and} \quad R^*(a) = \sum_{x: g(x)=a} \hat{\pi}(x) R(x)$$

provide (i) *cell-level* prioritization within \mathcal{X} and (ii) *technology-level* prioritization via R^* . By construction, R is monotone in each η_j and stable under \mathcal{N} .

Stability and calibration. (S1) *Normalization stability.* If $x \sim_{\mathcal{N}} x'$, then $\mathbf{f}(x) = \mathbf{f}(x')$ and $\eta_j(x) = \eta_j(x')$, hence $\pi_\theta(x) = \pi_\theta(x')$ and $R(x) = R(x')$.

(S2) *Threshold calibration.* The level-sets $\mathcal{A}_q(\tau)$ are nested: $\tau_1 \leq \tau_2 \Rightarrow \mathcal{A}_q(\tau_2) \subseteq \mathcal{A}_q(\tau_1)$. For any group a , the group-mass $\Gamma(a)$ is a right-continuous, piecewise-constant, nonincreasing function of τ on finite \mathcal{X} (i.e., changes only at finitely many breakpoints).

(S3) *Feasibility preservation.* If μ respects \mathcal{M} and \mathcal{N} (e.g., masking to Ω and canonical tokens), then $\hat{\pi}$ respects both by construction (projection onto $\mathcal{E} \cap \mathcal{M}$); in particular, no mass is assigned to infeasible axes or non-canonical tokens.

Explainability and audit. Let $\eta(x)$ be the feature vector and $R(x) = \sum_j \alpha_j \eta_j(x)$. The *principal reason* for a selection at x is $\arg \max_j \alpha_j \eta_j(x)$. Because η_j are counts/flags over (o, t, c, a) or local neighborhoods, explanations reduce to reporting sufficient statistics and typed support paths in G validating $x \in \mathcal{A}_q(\tau)$ (cf. Sec. 2). This yields compact, auditable rationales with no implementation dependence.

Optional diversity at presentation-time. If one wishes to select a small panel $\mathcal{S} \subset \mathcal{A}_q(\tau)$ of size K for display, define

$$F(A) = \sum_{x \in A} R(x) + \gamma \text{Div}(A), \quad |A| \leq K,$$

where Div is a monotone submodular diversity (e.g., facility-location on a similarity kernel over \mathcal{X}). Then the greedy algorithm achieves a $(1 - 1/e)$ approximation to $\max_{|A| \leq K} F(A)$ (Nemhauser et al., 1978; Lin and Bilmes, 2011), and DPP-based Div induces repulsion among near-duplicates (Kulesza and Taskar, 2012).

4. Gap Finder: Coverage Functionals and Interpretable Ranking

Axes, measure, and coverage. Let $\mathcal{X} = \mathcal{O} \times \mathcal{T} \times \mathcal{C} \times \mathcal{A}$ be the product space (organism, parent tissue, condition {Sf, Gnd}, assay). Normalization \mathcal{N} is idempotent and axis-consistent (cf. Sec. 3). Let $\kappa : \mathcal{X} \rightarrow \mathbb{N}$ be a multiplicity functional (count of distinct observed instances; $\kappa \equiv 0$ off Ω). Fix $m \in \mathbb{N}$; define the covered set and the *gap set*

$$\mathcal{C}\mathcal{O}\mathcal{V} = \{x : \kappa(x) \geq m\}, \quad \mathcal{G} = \mathcal{X} \setminus \mathcal{C}\mathcal{O}\mathcal{V}.$$

We endow \mathcal{X} with a product metric $d(x, y) = \sum_j w_j d_j(x_j, y_j)$ (discrete on categorical axes), and with the counting measure $\mu_\kappa(A) = \sum_{x \in A} \kappa(x)$.

Neighborhoods and feasibility. For $\rho > 0$ let $\mathcal{N}_\rho(x) = \{y \in \mathcal{X} : d(x, y) \leq \rho\}$ and define a normalized density

$$v_\rho(x) = \frac{1}{Z_\rho} \sum_{y \in \mathcal{N}_\rho(x)} \kappa(y), \quad Z_\rho = \sum_{y \in \mathcal{X}} \kappa(y).$$

Assay feasibility is encoded by $\varphi : \mathcal{A} \rightarrow [0, 1]$; set $\vartheta(x) = \varphi(a)$ for $x = \langle o, t, c, a \rangle$.

Signals (compact Greek block). For $x = \langle o, t, c, a \rangle \in \mathcal{G}$ define

$$\begin{aligned} \gamma(x) &= \min\left(1, \frac{\kappa(\langle o, t, \text{Gnd}, a \rangle)}{3}\right) \mathbb{1}\{c = \text{Sf}\}, \\ \mu(x) &= \mathbb{1}\{\exists a' \neq a : \kappa(\langle o, t, c, a' \rangle) \geq 1\}, \\ \phi(x) &= \mathbb{1}\{\exists c' \neq c : \kappa(\langle o, t, c', a \rangle) \geq 1\}, \\ \sigma(x) &= \mathbb{1}\{\exists o' \neq o : \kappa(\langle o', t, c, a \rangle) \geq 1\}, \\ v_\rho(x) &= \frac{1}{Z_\rho} \sum_{y \in \mathcal{N}_\rho(x)} \kappa(y), \quad \vartheta(x) = \varphi(a). \end{aligned}$$

(Abbreviations: Sf = Spaceflight; Gnd = Ground/Analog.)

Scoring functional and properties. Let $\alpha \in \mathbb{R}_{\geq 0}^6$ and $\beta \geq 0$. Define the compact score

$$\begin{aligned} S(x) &= \alpha_1 \gamma(x) + \alpha_2 \mu(x) + \alpha_3 \phi(x) \\ &\quad + \alpha_4 \sigma(x) + \alpha_5 v_\rho(x) + \alpha_6 \vartheta(x) - \beta \text{Red}(x). \end{aligned}$$

Here $\text{Red}(x)$ is a redundancy proxy (e.g., kernel similarity to previously selected gaps). Then:

- *Monotonicity in signals.* S is nondecreasing in each positive signal (by linear construction).
- *Normalization stability.* If $x \sim_{\mathcal{N}} x'$, then $\gamma, \mu, \phi, \sigma, v_\rho, \vartheta$ coincide, hence $S(x) = S(x')$.
- *Bounded sensitivity (local).* For x, y with $d(x, y) \leq \rho$, the variation of v_ρ is bounded by $\frac{\max_z \kappa(z)}{Z_\rho}$ times the discrepancy in neighborhood overlap; on finite domains this induces local stability of S .

Diversity-aware selection (budgeted). Given $K \in \mathbb{N}$, select a panel $A \subseteq \mathcal{G}$, $|A| \leq K$, by

$$F(A) = \sum_{x \in A} S(x) + \gamma \text{Div}(A),$$

where Div promotes dispersion:

$$\text{Facility-location: } \text{Div}(A) = \sum_{y \in \mathcal{G}} \max_{x \in A} k(x, y),$$

$$\text{DPP (log-MAP): } \text{Div}(A) = \log \det(K_A),$$

with a PSD kernel $K = (k(x, y))$ over \mathcal{G} . If Div is monotone submodular, the cardinality-constrained greedy achieves $(1 - 1/e)$ (Nemhauser et al., 1978; Krause and Golovin, 2014). For DPPs, greedy-MAP selections encourage explicit repulsion and diversity (Kulesza and Taskar, 2012).

Connectivity repair (optional). If connectivity is desired in projections of A (e.g., by tissue), apply a minimal post-augmentation via a Steiner-style repair on the co-occurrence graph over \mathcal{X} ; this keeps the score close to $F(A)$ in practice, and the repair can leverage approximation algorithms for Steiner trees (Vazirani, 2003; Robins and Zelikovsky, 2000).

Explainability as sensitivity decomposition. Let $s(x) = (\gamma, \mu, \phi, \sigma, v_\rho, \vartheta)$ and $S(x) = \langle \alpha, s(x) \rangle - \beta \text{Red}(x)$. *Principal reason:* $\arg \max_j \alpha_j s_j(x)$. *Detail:* report $s(x)$ and, when applicable, underlying statistics (e.g., counts in $\mathcal{N}_\rho(x)$). Auditability is completed with typed paths in G that justify the signals (cf. Sec. 2).

Calibration and scope. Let $\mathcal{X}_{\text{sc}} \subseteq \mathcal{X}$ be the *scope* induced by filters (or the observed support). Then: (i) κ and v_ρ are computed over \mathcal{X}_{sc} ; (ii) S is invariant to extensions outside \mathcal{X}_{sc} that do not modify κ on $\mathcal{N}_\rho(x)$; (iii) increasing m shrinks $\mathcal{E} \mathcal{O} \mathcal{V}$ and expands \mathcal{G} , with S stable when $\gamma, \mu, \phi, \sigma$ remain unchanged.

Relation to Assay Finder. Assay Finder defines $\hat{\pi}$ over \mathcal{X} (Sec. 3). An optional compatible prior rescales S as $\tilde{S}(x) = \hat{\pi}(x) S(x)$, focusing the panel on conceptually relevant cells without altering monotonicity or normalization stability.

5. Mission Builder: Multiobjective Planning with Evidence-Constrained Synthesis

Objects and evidence. Let \mathcal{O} be a finite set of mission objectives, \mathcal{P} a finite set of phases (e.g., PDR \rightarrow CDR \rightarrow Ops), and \mathcal{R} a resource index (mass, fuel, crew, power, ...). A *mission architecture* is a tuple

$$\mathbf{m} = (\Theta, \Pi, \mathcal{T}, \rho, \xi),$$

where $\Theta \subseteq \mathcal{O}$ (selected objectives), Π is a phasewise policy, \mathcal{T} is a precedence-respecting timeline, $\rho \in \mathbb{R}_{\geq 0}^{|\mathcal{R}|}$ is the aggregated resource vector, and $\xi \in [0, 1]$ is an aggregated risk level. GraphRAG (Sec. 2) contributes *typed evidence* $E = \{(c_i, w_i)\}_{i=1}^M$ with weights $w_i \in [0, 1]$ that support Θ and key decisions in Π via typed shortest paths in G .

Timeline and precedence. Let (\mathcal{P}, \preceq) be a precedence DAG. A *schedule* is a map $t : \mathcal{P} \rightarrow \mathbb{R}_{\geq 0}$ such that

$$p \preceq p' \Rightarrow t(p) + d(p) \leq t(p'),$$

where $d(p)$ is the canonical duration of p . Define aggregate slack $\Delta = \sum_{p \in \mathcal{P}} \max\{0, t_{\max}(p) - t(p) - d(p)\}$ (Pinedo, 2016).

Resources and feasibility. Let $r_p \in \mathbb{R}_{\geq 0}^{|\mathcal{R}|}$ be phase consumption and $r = \sum_p r_p$ the total; let $\bar{r} \in \mathbb{R}_{\geq 0}^{|\mathcal{R}|}$ be limits.

$$\text{Feasibility: } r \leq \bar{r} \quad (\text{componentwise}).$$

For each objective $\theta \in \Theta$, specify minimal requirements $r(\theta)$ and a set of enabling phases $\mathcal{P}(\theta)$ (as in RCPSP-style feasibility) (Brucker et al., 1999; Herroelen and Leus, 2005).

Risk model (coherent). Let Ω be a finite scenario set with probability \mathbb{P} . Each $p \in \mathcal{P}$ induces a nonnegative loss $L_p(\omega)$;

total $L(\omega) = \sum_p L_p(\omega)$. Define risk by a coherent risk measure, e.g., CVaR $_\alpha$:

$$\xi = \text{CVaR}_\alpha(L) = \min_{\tau \in \mathbb{R}} \tau + \frac{1}{1-\alpha} \mathbb{E}[(L - \tau)_+],$$

with coherence per Artzner et al. (1999) and convex representation per Rockafellar and Uryasev (2000, 2002). For stochastic planning background, see Shapiro et al. (2009).

Utility and evidence constraints. Let $u : \Theta \rightarrow \mathbb{R}_{\geq 0}$ denote objective utility and $g : \mathcal{P} \rightarrow \mathbb{R}_{\geq 0}$ denote gains for critical phases. The structural utility is

$$U(\Theta, \Pi) = \sum_{\theta \in \Theta} u(\theta) + \sum_{p \in \mathcal{P}} g(p) \mathbb{1}_{\{\Pi \text{ executes } p\}}.$$

Evidence \mathbf{E} imposes *consistency*: there exists a set $\mathcal{Q} \subseteq \mathcal{P}(G)$ of typed paths in G such that

$$\forall \theta \in \Theta \exists q \in \mathcal{Q} : q \models \theta, \quad \sum_{(c_i, w_i) \in \mathbf{E} : c_i \in q} w_i \geq \tau_{\text{ev}},$$

and analogously for decisions in Π (justification of instruments/resources). This enforces architecture-level provenance-closure (cf. Sec. 2).

Multiobjective synthesis (scalarized). With weights $\lambda = (\lambda_U, \lambda_\Delta, \lambda_\rho, \lambda_\xi)$, define

$$J(\mathbf{m}) = \lambda_U U(\Theta, \Pi) + \lambda_\Delta \Delta - \lambda_\rho \|r\|_1 - \lambda_\xi \xi.$$

Mission Builder solves

$$\max_{\mathbf{m}} J(\mathbf{m}) \quad \text{s.t.} \quad \begin{cases} \text{precedences and durations,} \\ r \leq \bar{r}, \quad \theta \in \Theta \Rightarrow \mathcal{P}(\theta) \subseteq \Pi, \\ \text{consistency with } \mathbf{E}, \quad \xi \leq \bar{\xi}. \end{cases}$$

Pareto fronts can be traced by varying λ (weighted-sum scalarization) or via ε -constraint formulations (Miettinen, 1999; Haimes et al., 1971). (Weighted sums recover supported efficient points under standard convexity assumptions.)

Context packing for rationale (compact). Let $\text{ctx} \subseteq \mathbf{E}$ be a *budgeted* evidence portfolio to accompany the architecture. Select ctx by

$$\max_{|\text{ctx}| \leq B} \text{DPP}(\text{ctx}) + \text{Cov}(\text{ctx}; \Theta, \Pi),$$

where DPP induces repulsion and Cov is a coverage functional over objectives/phases (Kulesza and Taskar, 2012; Nemhauser et al., 1978; Krause and Golovin, 2014). This yields compact, nonredundant rationales within budget B .

Phase policies (Markovian abstraction). Let \mathcal{S} be a high-level state space with $s_{t+1} = F(s_t, a_t, \omega_t)$ and $a_t \in \mathcal{A}_t$ (phase decisions). A policy $\Pi = (\pi_t)$ is Markovian if $a_t = \pi_t(s_t)$. For finite horizons, a stochastic formulation is

$$\max_{\Pi, t(\cdot)} \mathbb{E} \left[\sum_t r_t(s_t, a_t) - \lambda_\xi L(\omega) \right]$$

s.t. precedences, $r \leq \bar{r}$, consistency with \mathbf{E} .

in the spirit of constrained MDPs and risk-aware planning (Puterman, 1994; Shapiro et al., 2009). The realized policy is “compiled” to \mathbf{m} by extracting achieved Θ , \mathcal{T} , and ρ .

Feasibility cuts and repair. If a candidate violates evidence or resources, add cuts of the form

$$\sum_{p \in \mathcal{P}(\theta)} z_p \geq 1 \quad (\theta \in \Theta), \quad \sum_p a_{pr} z_p \leq \bar{r}_r \quad (r \in \mathcal{R}),$$

with $z_p \in \{0, 1\}$ execution indicators (valid inequalities/Benders-style logic cuts) (Geoffrion, 1972). Precedence repairs can be posed as minimal augmentations on the temporal DAG (using Steiner-style heuristics); this is a pragmatic post-step leveraging approximation algorithms, rather than introducing new guarantees (Vazirani, 2003; Robins and Zelikovskiy, 2000).

Compact invariants (column-friendly).

$$(I) \forall \theta \in \Theta \exists q \in \mathcal{Q} : q \models \theta \quad (\text{provenance})$$

$$(II) r \leq \bar{r}, \quad \xi \leq \bar{\xi}, \quad t \text{ respects } (\mathcal{P}, \preceq) \quad (\text{feasibility})$$

$$(III) \max_{\text{ctx}} \{\text{DPP} + \text{Cov}\} \text{ s.t. } |\text{ctx}| \leq B \quad (\text{rationale diversity})$$

(II) separates resources, risk, and time; (III) aligns explanation with diversity/coverage principles.

Manager-oriented views (scores). Define a *mission score*

$$\mathcal{S}(\mathbf{m}) = \beta_1 U + \beta_2 \Delta - \beta_3 \|r\|_1 - \beta_4 \xi,$$

and an optional *gap-closure gain* $\sum_{x \in \mathcal{G}} \tilde{S}(x) \mathbb{1}_{\{\text{closed by } \mathbf{m}\}}$ (cf. Sec. 4). Both are traceable: \mathcal{S} decomposes by objectives/phases, and each term is justified by ctx and typed paths in G .

References

- Amari, S.-i., Nagaoka, H., 2000. Methods of Information Geometry. AMS & Oxford University Press.
- Andersen, R., Chung, F., Lang, K., 2006. Local graph partitioning using pagerank vectors. In: FOCS.
- Artzner, P., Delbaen, F., Eber, J.-M., Heath, D., 1999. Coherent measures of risk. Mathematical Finance 9 (3), 203–228. DOI: 10.1111/1467-9965.00068
- Blondel, V. D., Guillaume, J.-L., Lambiotte, R., Lefebvre, E., 2008. Fast unfolding of communities in large networks. Journal of Statistical Mechanics.
- Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., Yakhnenko, O., 2013. Translating embeddings for modeling multi-relational data. In: NeurIPS.
- Brown, L. D., 1986. Fundamentals of Statistical Exponential Families with Applications in Statistical Decision Theory. IMS.
- Brucker, P., Drexel, A., Möhring, R., Neumann, K., Pesch, E., 1999. Resource-constrained project scheduling: Notation, classification, models, and methods. European Journal of Operational Research 112 (1), 3–41. DOI: 10.1016/S0377-2217(98)00204-5
- Chung, F. R. K., 1997. Spectral Graph Theory. American Mathematical Society.
- Csiszár, I., 1975. I-divergence geometry of probability distributions and minimization problems. The Annals of Probability 3 (1), 146–158. DOI: 10.1214/aop/1176996454
- Edge, D., Trinh, H., Cheng, N., Bradley, J., Chao, A., Mody, A., Truitt, S., Larson, J., 2024. From local to global: A graph rag approach to query-focused summarization. arXiv:2404.16130.
- Gao, Y., et al., 2023. Retrieval-augmented generation for large language models: A survey. arXiv:2312.10997.
- Geoffrion, A. M., 1972. Generalized benders decomposition. Journal of Optimization Theory and Applications 10 (4), 237–260. DOI: 10.1007/BF00934810

- Haimes, Y. Y., Lasdon, L. S., Wismer, D. A., 1971. On a bicriterion formulation of the problems of integrated system identification and system optimization. *IEEE Transactions on Systems, Man, and Cybernetics SMC-1* (3), 296–297.
DOI: 10.1109/TSMC.1971.4308298
- Herroelen, W., Leus, R., 2005. Project scheduling under uncertainty: Survey and research potentials. *European Journal of Operational Research* 165 (2), 289–306.
DOI: 10.1016/j.ejor.2004.04.002
- Hogan, A., Blomqvist, E., Cochez, M., et al., 2021. Knowledge graphs. *ACM Computing Surveys* 54 (4), 1–37.
DOI: 10.1145/3447772
- Jaynes, E. T., 1957. Information theory and statistical mechanics. *Physical Review* 106 (4), 620–630.
DOI: 10.1103/PhysRev.106.620
- Krause, A., Golovin, D., 2014. Submodular function maximization. In: *Tractability: Practical Approaches to Hard Problems*. Cambridge University Press.
- Kulesza, A., Taskar, B., 2012. *Determinantal Point Processes for Machine Learning*. Now Publishers.
- Larson, J., Truitt, S., February 2024. Graphrag: Unlocking llm discovery on narrative private data.
- Lin, H., Bilmes, J., 2011. Class-subset selection via submodular augmentation. In: *Proc. NIPS*.
- Manning, C. D., Raghavan, P., Schütze, H., 2008. *Introduction to Information Retrieval*. Cambridge University Press.
- Microsoft Research, 2024. Project graphrag. Project page.
- Miettinen, K., 1999. *Nonlinear Multiobjective Optimization*. Springer.
DOI: 10.1007/978-1-4615-5563-6
- Mungall, C. J., et al., 2012. Uberon: an integrative multi-species anatomy ontology. *Genome Biology* 13 (1), R5.
- Nemhauser, G. L., Wolsey, L. A., Fisher, M. L., 1978. An analysis of approximations for maximizing submodular set functions. *Mathematical Programming* 14, 265–294.
DOI: 10.1007/BF01588971
- Newman, M. E. J., 2006. Modularity and community structure in networks. *PNAS* 103 (23), 8577–8582.
- Ng, A. Y., Jordan, M. I., Weiss, Y., 2001. On spectral clustering: Analysis and an algorithm. In: *NeurIPS*.
- Peng, B., et al., 2024. Graph retrieval-augmented generation: A survey. *arXiv:2408.08921*.
- Pinedo, M. L., 2016. *Scheduling: Theory, Algorithms, and Systems*, 5th Edition. Springer.
DOI: 10.1007/978-3-319-26580-3
- Puterman, M. L., 1994. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley.
- Robins, G., Zelikovsky, A., 2000. Improved steiner tree approximation in graphs. *SIAM Journal on Discrete Mathematics* 19 (1), 122–134.
DOI: 10.1137/S0895480100375002
- Rockafellar, R. T., Uryasev, S., 2000. Optimization of conditional value-at-risk. *Journal of Risk* 2 (3), 21–41.
- Rockafellar, R. T., Uryasev, S., 2002. Conditional value-at-risk for general loss distributions. *Journal of Banking & Finance* 26 (7), 1443–1471.
DOI: 10.1016/S0378-4266(02)00271-6
- Shapiro, A., Dentcheva, D., Ruszczyński, A., 2009. *Lectures on Stochastic Programming: Modeling and Theory*. SIAM.
DOI: 10.1137/1.9780898718751
- Shi, J., Malik, J., 2000. Normalized cuts and image segmentation. *TPAMI* 22 (8), 888–905.
- Smith, B., et al., 2007. The obo foundry: Coordinated evolution of ontologies to support biomedical data integration. *Nature Biotechnology* 25 (11), 1251–1255.
- Sun, Z., Deng, Z.-H., Nie, J.-Y., Tang, J., 2019. Rotate: Knowledge graph embedding by relational rotation in complex space. In: *ICLR*.
- Tong, H., Faloutsos, C., Pan, J.-Y., 2006. Fast random walk with restart and its applications. In: *ICDM*. pp. 613–622.
- Trouillon, T., Welbl, J., Riedel, S., Gaussier, É., Bouchard, G., 2016. Complex embeddings for simple link prediction. In: *ICML*.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., Polosukhin, I., 2017. Attention is all you need. In: *Advances in Neural Information Processing Systems (NeurIPS)*.
- Vazirani, V. V., 2003. *Approximation Algorithms*. Springer.
DOI: 10.1007/978-3-662-04565-7
- Wainwright, M. J., Jordan, M. I., 2008. Graphical models, exponential families, and variational inference. In: *Foundations and Trends in Machine Learning*. Vol. 1. pp. 1–305.
DOI: 10.1561/22000000001
- Yu, H., et al., 2024. Evaluation of retrieval-augmented generation: A survey. *arXiv:2405.07437*.
- Zhang, Q., et al., 2025. A survey of graph retrieval-augmented generation for llms. *OpenReview*.