

# An Extensible Model-Based Mediator System with Domain Maps

Amarnath Gupta\*

Bertram Ludäscher\*

Maryann E. Martone<sup>‡</sup>

\*San Diego Supercomputer Center, UCSD {gupta,ludaesch}@sdsc.edu

<sup>‡</sup>Department of Neurosciences, UCSD mmartone@ucsd.edu

## 1 Background

Mediator systems federate and integrate data from disparate sources in order to elicit information that the individual sources cannot provide independently. The standard mediator architecture employs wrappers that translate heterogeneous source data into a common (often semistructured) data model like XML. A “mediation engineer” provides an integrated view definition (IVD) on the wrapped XML sources. In such a system, an IVD is ideally expressed in a declarative query language for XML or semistructured data. When developing the IVD, an XML query language provides the mediation engineer only with a *tree-structured* model of the source, *i.e.*, the names and possible nesting structure of XML elements as defined by an XML DTD, but gives no hint on semantic relationships, class structures, not to mention application domain specific constraints. Indeed, as shown in [4, 1, 3], mediation should be lifted to the conceptual level when mediating across complex sources whose data comes from seemingly disjoint “worlds”, *e.g.*, two neuroscience labs creating information on neurotransmitters and protein distributions, respectively.<sup>1</sup>

To this end, we present a mediator prototype system whose main novel features are: (i) mediated views are defined *and executed* at the level of *conceptual models* (CMs) rather than at the usual structural level, (ii) *domain maps* (DMs) – labeled graphs of concepts and relationships with a formal logic semantics – are used to bridge the semantic gap between source data from “multiple worlds”, and (iii) a *plug-in mechanism* for CMs and DMs is provided which allows the mediator system to be easily extended when new formalisms for CMs and DMs are used by sources. We illustrate these features using an example from a complex neuroscience mediation problem.

For details of model-based mediation with domain maps, including their formal semantics, see [4].

<sup>1</sup>See [senselab.med.yale.edu](http://senselab.med.yale.edu) and [www-ncmir.ucsd.edu](http://www-ncmir.ucsd.edu).

## 2 Model-Based Mediation

Figure 1 depicts our system architecture for model-based mediation: Differences in the sources’ data models are resolved by wrappers that translate the raw data into a common generic data format (XML). Current mediator systems directly define the integrated views on the wrapped XML sources using an XML query language. We extend this architecture by lifting exported source data to the semantically richer level of conceptual models with domain knowledge. Thus, the integrated view definition IVD at the mediator is aware of class hierarchies, object structure, properties of relationships (inclusion dependencies, cardinalities, ...), and in particular *domain specific constraints* of sources. Consequently, the mediator’s view definition language in this architecture must not only act as query language for semistructured data, but also for conceptual models including the definition of complex schema and instance level transformations and checking of logical constraints.

### 2.1 Generic Conceptual Model (GCM)

To facilitate extensibility, we use GCM, a generic conceptual model, at the mediator level. Like RDF, GCM is a minimalist object-oriented model that allows specification of objects at the schema and instance level (*e.g.*, `method(classA,meth,classB)` and `methodinst(oidA,meth,oidB)`), similar for relations, and – most importantly – a *rule-based extension mechanism* for axiomatizing additional CM constructs and constraints. The formal model of GCM is a fragment of F-logic (short: FL) [2] with well-founded negation semantics. In this way, GCM is *universal* for CMs since all first-order constraints (cardinality constraints, range constraints, inclusion dependencies, *etc.*) and the usual inductive constraints (*e.g.*, transitivity of the class hierarchy) are expressible in the GCM formalism. The choice of a FL for our GCM is partly for convenience, since FL already includes all required GCM features and we thus get a GCM for-



respond to anatomical entities of the brain, edges correspond to relationships like *is-a* and *has-a*. Moreover, different shades indicate absence or (in)direct presence of data. Like CMs, DMs may have additional rules, in this case, that *has-a* needs to be closed wrt. the transitive *is-a* relation; the result is shown at the bottom of Figure 2.

## 2.4 Query Processing

At runtime, a wrapped source  $S$  joins the mediation by registering its conceptual model  $CM(S)$  with the mediator  $M$ . This requires that  $S$  sends the mediator descriptions of the exported class schemas, relationship schemas, and semantic rules (ultimately expressed in FL) that are evaluable at the mediator (either using GCM, or any CM formalism for which a plug-in is available). The exported objects of a  $CM(S)$  can have special *context* attributes that provide the “semantic coordinates” of the data in the mediator’s *domain map*  $DM(M)$ . In particular, the *context* attributes can create new concept classes at  $M$  as a result of the source’s registration process.

At the mediator, the user query is executed against the IVD. For example, a plan for the user query “*What is the distribution of those calcium-binding proteins that are found in neurons that receive signals from parallel fibers in rat brains?*” involves the following steps: (1) *push selections* (*‘rat’*, *‘parallel fiber’*) to the SENSELAB source and *get bindings* for neuron/compartment pairs  $X$  and  $Y$ ; (2) using the domain map  $DM(M)$ , *select sources* that have data anchored at  $X, Y$  from step (1) – in our case, NCMIR; (3) *push selections* given by the  $X, Y$  locations to NCMIR, and *retrieve* only proteins  $P$  that are found in  $X, Y$ ; (4) based on the *least upper bound* of locations in the domain map, *compute the view* *protein\_distribution* at the mediator (this involves a *downward closure* along the *has-a-star* relation).

The last two operations filter out a segment in the domain map as the “region of correspondence” between the two information sources, and demonstrate how graph operations on the domain map can be actively used to compute conceptual mappings between sources (cf. Figure 3).

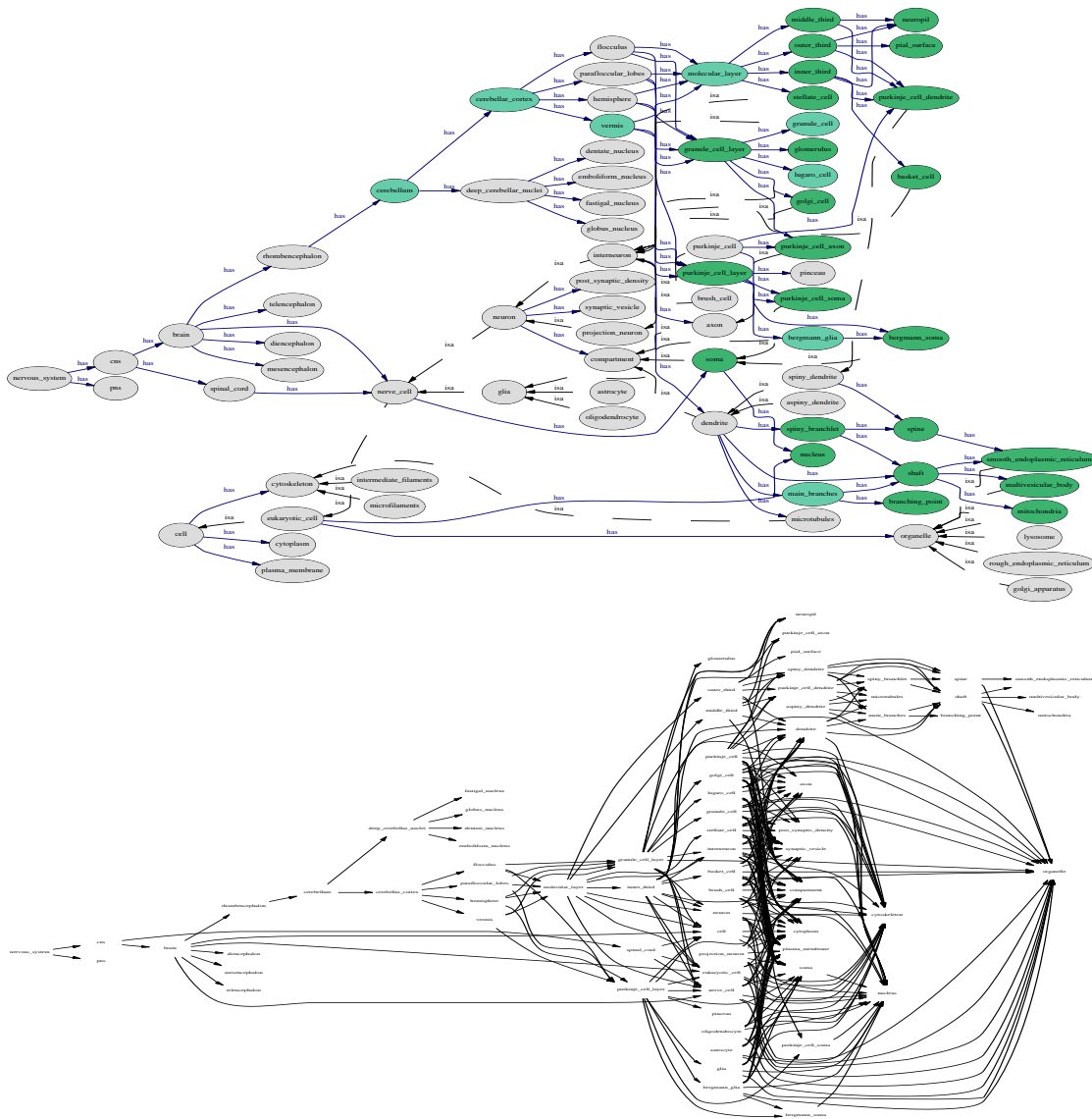
For the full

## References

[1] A. Gupta, B. Ludäscher, and M. E. Martone. Knowledge-Based Integration of Neuro-

science Data Sources. *12th Intl. Conf. on Scientific and Statistical Database Management (SSDBM)*, July 2000.

- [2] M. Kifer, G. Lausen, and J. Wu. Logical Foundations of Object-Oriented and Frame-Based Languages. *Journal of the ACM*, 42(4):741–843, July 1995.
- [3] B. Ludäscher, A. Gupta, and M. E. Martone. Model-Based Information Integration in a Neuroscience Mediator System. In *26th Conf. on Very Large Data Bases (VLDB)*, Cairo, Egypt, 2000. system demonstration.
- [4] B. Ludäscher, A. Gupta, and M. E. Martone. Model-Based Mediation with Domain Maps. In *Intl. Conf. on Data Engineering (ICDE)*, [www.sdsc.edu/~ludaesch/Paper/icde01.html](http://www.sdsc.edu/~ludaesch/Paper/icde01.html), 2001.
- [5] J. Suzuki and Y. Yamamoto. Making UML Models Interoperable with UXF. In *Intl. Workshop <<UML’98>>: Beyond the Notation*, LNCS 1618, 1998.
- [6] G. Yang and M. Kifer. FLORA: Implementing an Efficient DOOD System Using a Tabling Logic Engine. In *6th International Conference on Rules and Objects in Databases (DOOD)*, 2000.



**Figure 2.** ANATOM domain map: is-a  $\cup$  has-a (above), and has-a-star, the deductive-closure of is-a and has-a (below)

The screenshot displays the UCSD Neuro-Centroid-DB Demo interface. On the left is a navigation menu with categories like HOME, SLIDES, MEDIATED, CONCEPTUAL, SCHEMA, ANATOMY, and others. The main content area is divided into three sections: 1. Interactive PROTEIN LOCALIZATION (UCSD/NCMIR data only), 2. Look (Show Slides), and 3. Inter (lookin' How does Receptor?).

The 'Protein Localization Results' section features a table with the following data:

PROTEIN	ISOFORM	SPECIES	IMAGE ID
Ryanodine Receptor	?	rat	?

Below the table are two image viewers: 'Image Objects' and 'Surface Objects'. The 'Image Objects' viewer shows a grayscale image of a brain section with a red box highlighting a region. The 'Surface Objects' viewer shows a 3D surface reconstruction of the same region. A table below these viewers provides quantitative data for the objects:

Object	Perimeter	Area	Centroid X	Centroid Y	Centroid Z
dd_sdr_rfrag	79.7859	104.7799	414.8138	122.6698	0.0000
neuro_t...	21.8700	23.2600	442.1599	52.6800	0.0000
...	39.7700	22.1200	362.8900	239.1599	0.0000
...	44.2000	39.2500	349.2999	253.9199	0.0000
...	49.6400	47.8799	366.3099	329.0199	0.0000
...	45.9300	45.8100	427.9399	55.2100	0.0000
...	32.6299	35.0900	428.7999	68.4933	0.0000
...	45.8000	55.7900	413.4999	76.9800	0.0000
...	49.2400	59.4600	389.0297	230.8398	0.0000
...	...	...	...	...	...

On the right side, a subgraph of anatomical data is shown, starting with 'nervous\_system' and branching into 'neu', 'enu', 'brain', 'thalamencephalon', 'cerebrum', 'cerebellar\_cortex', 'vermis', 'molecular\_layer', 'inner\_third', 'middle\_third', 'outer\_third', 'purkinje\_cell\_dendrite', 'spiny\_branchelet', 'mean\_branchelet', and 'branching\_point'. Two diamond-shaped nodes at the bottom indicate '[108 results]' and '[36 results]'.

Figure 3. Snapshot of the mediator prototype; *background left*: mediator shell for issuing ad-hoc queries against CM(M); *background right*: generated subgraph of ANATOM having the requested result data; clicking on a (diamond) result node retrieves the actual result data (*foreground center*)