# Domain-guided construction of semantic representations for model-based interpretation

*Cleo Condoravdi*
CSLI, Stanford University

*Richard Waldinger*
SRI

Joint work with D. Bobrow, K. Richardson, A. Das

# Quadri Project Team

Funding: National Institute of Health (NIH)

- PARC
  - **Danny Bobrow**
  - **Cleo Condoravdi** (now at Stanford University)
  - **Kyle Richardson** (now at University of Stuttgart)
- SRI International
  - **Richard Waldinger** Artificial Intelligence Center
- Stanford University
  - **Amar Das** Biomedical Informatics Research
  - **Bob Shafer** Stanford HIV Database Curator
  - **Soo-Yon Rhee** Stanford HIV Database Curator

# Textual Inference Task

Does premise *P* lead to conclusion *C*?
Does text *T* support the hypothesis *H*?
Does text *T* answer the question *H*?
    *... without any additional assumptions*

**P***: Every explorer failed to get to the South Pole.*
**C***: No experienced explorer reached the South Pole.*
  Yes

**P***: Ed has been living in Athens for 3 years.*
  *Mary visited Athens in the last 2 years.*
**C***: Mary visited Athens while Ed lived in Athens.*
  Yes

# Inference Task

Does a given specifications of the world $D$ support the statement $S$?

Is statement $S$ true relative to a state of the world as specified by $D$?

What is the answer to the question $Q$ relative to a dataset/ database $D$?

*Which rivers flow through the states that border California*?

# Geobase

A small database of information about United States geography with about 800 facts, represented as Prolog assertions

States - their capitals, populations, areas, population densities, major cities, rivers and the bordering states

Cities - their populations and the states they are in

Rivers - their lengths and the states through which they flow

Mountains - their heights and the states they are in

# Inference Task

What is the answer to question *Q* relative to a dataset/database *D*?

http://www.cs.utexas.edu/users/ml/geo-demo.html

**Geoquery:**

*Which rivers flow through the states that border California*?

**CHILL:**

**[colorado,columbia,gila,snake]**

Formal Language Query:

answer(_74,
(river(_74),traverse(_74,_75),state(_75),next_to(_75,_76),const(_76,stateid(california))))

X borders Y ➡ X next_to Y
X flows through Y ➡ X traverse Y

# Inference Task

What is the answer to question *Q* relative to a dataset/database *D*?

**Geoquery:**
*How many states does the Mississipi run through***?**
**CHILL:**
**[10]**

Formal Language Query:

answer(_86,count(_87,
(state(_87),const(_88,riverid(mississippi)),traverse(_88,_87)),_86))

# Inference Task

What is the answer to question *Q* relative to a dataset/database *D*?

http://www.cs.utexas.edu/users/ml/geo-demo.html

**Geoquery:**

*Does California have at least 2 rivers*?

**CHILL:**

**[mississippi]**

Formal Language Query:

answer(_82,(const(_83,stateid(california)),smallest(_83,river(_82))))

at least 2 rivers ➡ cardinality of rivers such that …
at least 2 rivers ➡ smallest river

# Database table structure: temporally bound treatments

## Regimen Table

| Field | DB Type |
|---|---|
| Regimen Id | Key Id |
| Patient Id | Id |
| Start Date | String (D/M/Y) |
| End Date | String (D/M/Y) |
| Drug Set | String (D1+D2+D3) |

*What regimens include drug AZT?*
*What patients had a regimen with at least 2 PIs?*
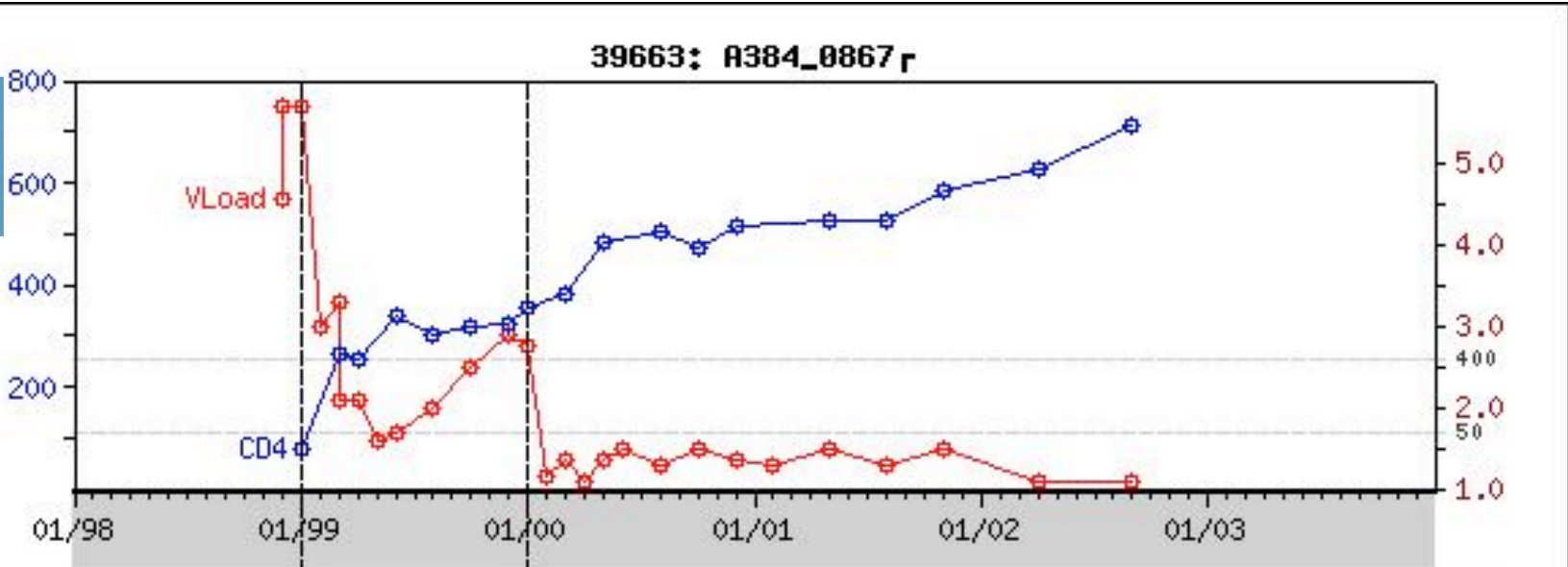*What patients had a regimen with EFV for more than 24 weeks?*

# HIV drug resistance

- HIV has complex treatment patterns
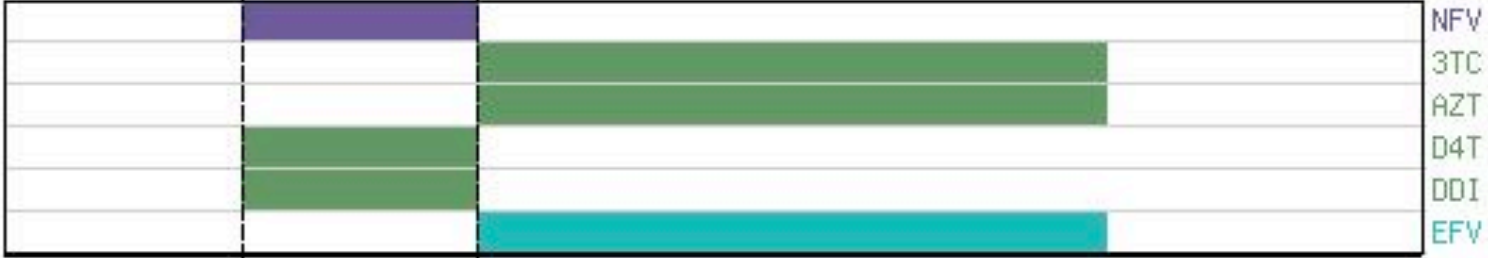- Drug-resistant mutations are a major obstacle to the success of treatment

- Stanford has useful databases in this domain
- Anonymized patient records
- Summaries of clinical trials
- Ontologies of drugs, treatments, terms

# HIV Drug Resistance

# HIV Drug Resistance

# Database table structure: temporally bound treatments

## Regimen Table

| Field | DB Type |
|-------|---------|
| Regimen Id | Key Id |
| Patient Id | Id |
| Start Date | String (D/M/Y) |
| End Date | String (D/M/Y) |
| Drug Set | String (D1+D2+D3) |

*What regimens include drug AZT?*
*What patients had a regimen with at least 2 PIs?*
*What patients had a regimen with EFV for more than 24 weeks?*

# Virtual tables support higher level queries

**TCE (treatment change episode) Table**

| Field | DB Type |
|---|---|
| TCE Id | Key Id |
| Patient Id | Id |
| Failing Reg. | Id |
| Salvage Reg. | Id |
| Start Date | String (D/M/Y) |
| End Date | String (D/M/Y) |
| Baseline Duration | Number |

*What TCEs have a genotype of M184V during the failing regimen?*

# Motivation for NL Interface to databases
# How can I see what is in those databases?



What patients on Atripla exhibited a high viral load?

Stanford HIV clinical data

# What makes it difficult to access?

What patients on Atripla exhibited a high viral load?

What are the databases that are available?

What is their structure?

How do I get information out of them?

Multiple Databases

# Quadri: Intelligent Question Answering in the HIV Domain

*Natural Language Processing*

*Subject Domain Reasoning*

*Question Answering about Drug Resistance Information*

*Temporal Representation & Reasoning*

*Clinical Databases*

# Quadri: simplifies access in HIV domain
## Customizing general NL and Reasoning Systems



What patients on Atripla exhibited a high viral load?

PARC's Bridge

SRI's Snark

Stanford's HIV Databases + Other Resources

# Transformations in processing a query


Language Processing

- Text query
- Dependency parse
- Abstract KR
- Flat logical form (LF)
  with domain-specific relations


Logic Processing

- Translation to nested LF
- Feedback to user
- Prove the theorem
  domain theory + DB facts
- Display the answer

# Quadri architecture

# Sample questions

- *What mutations were found in patients after they failed AZT?*

- *Find all patients who had a high viral load on a regimen with EFV after 24 weeks.*

- *Find patients who were on Atripla for at least 12 weeks. They failed that regimen. They were then switched to a new regimen.*

-

# Axiomatic Subject-Domain Theory

- A domain-specific knowledge base where knowledge is expressed as *axioms*
- Higher level abstraction of the contents of the databases
  - Basic domain relations for which there is a correspondence in the databases, e.g. patient, patient-has-regimen
  - Derived domain relations, e.g. failing-regimen, AZT-naive
  - Translate qualitative specifications into quantitative specifications
  - Temporal axioms
  - Axioms relating regimens and their time spans

# HIV Domain

# Language Use Model

*English:*
*Patient = {'patient', ..}*
*Drug = {'epivir', 'norvir', …}*
*Regimen = {'regimen', 'treatment',..}*
*medicalTest = {'viral_load', 'enotype',..}*

**DATABASE**

**Patient <PatientID**, Region,...>
**RNA** <PATIENTID,**RNA_DATE,**
**VIRAL_LOAD_VOL>**
**Regimen**<PatientID,**Start_Date,Drug_List, ..>**
 ……

*Sorts = {Patient, medical_test,*
*Drug, Regimen, ….}*

*Relations:*
*(patient, regimen, patient-has-regimen)*
*(regimen, drug, regiman-has-drug)*
*(patient, medical_test, patient-has-test)*
*(medical_test, value, MT-has-value)*
*….*

# Semantic link to databases

- Link symbol in axiomatic theory with database(s)
- Axiomatic "advertisements" describe content of database
- The ground formulas of the theory are the relations in the database(s)
- Procedural attachments convert from date stamps in the database to time intervals
- Database invoked as proof search is underway

# Semantic types in the language

**Regimen Table**

| Field | DB Type | Semantic Type |
|---|---|---|
| Field | DB Type | Regimen |
| Regimen Id | Key Id | |
| Patient Id | Id | Patient |
| Start Date | String (D/M/Y) | Time Point |
| End Date | String (D/M/Y) | Time Point |
| Drug Set | String (D1+D2+D3) | Drug |

*What regimens include drug AZT?*

*What patients had a regimen with at least 2 PIs?*

*What patients had a regimen with EFV for more than 24 weeks?*

# Reasoning needed to interpret query
## *Find patients who had a **high viral load** after 24 weeks on a regimen with **Atripla**.*

Interpret qualitative terms wrt numbers

  **high** viral load *means* viral_load > **1000**

Expand **Atripla** wrt standard drugs

  EFV/FTC/TDF

efavirenz,emtricitabine, and tenofovir disoproxil fumarate

# Example Axiom

(failing-regimen-for-patient ?regimen ?patient
 ?time-point ?viral-load)

⟺

(and (patient-on-regimen ?patient ?regimen)
 (has-test viral-load ?patient ?time-point ?viral-load)
 (near-end ?time-point ?regimen)
 (viral-load-has-level ?viral-load high))

*A failing regimen for a patient is one in which the patient has a high viral load near the end of the regimen*

# Example Axiom

(near-end ?time-point ?time-interval)

⇔

(and

(within-pi ?time-point ?time-interval)
(=< (* 4 (minus-time (finish-time ?time-interval) ?time-point))
(duration ?time-interval))

*A time-point is near the end of a time-interval if it is in the 4th quarter of the interval* (can be changed)

# Quantitative reasoning about time

*Find patients who had a high viral load after 24 weeks on the regimen with Atripla*



24 weeks = 164 days

t

t' = date of test

Start of Regimen

Viral_load = high

Regimen with Atripla

…..

# Temporal Reasoning

- Reasoning about time points and intervals (Allen calculus)
- Date and time computations
- Durations
- Unit conversion

# BRIDGE system for language analysis

- BRIDGE is a broad coverage, general purpose natural language processing system.

- Stages of processing
  - Parsing
  - Abstract Knowledge Representation

- Bridge preserves ambiguities, marking local choices (packing).

- Customization
  - Task – Building Logical Forms
  - Domain – Recognizing HIV relations

# Parsing produces functional structures

*Find patients who were on Atripla.*



```
F-structure chart

"find patients who were on Atripla."

    ┌PRED    'find<[11-SUBJ:null_pro], [44:patient]>'                        ┐
    │SUBJ    ┌PRED 'null_pro']                                              │
    │        ┌PRED        'patient'                                         │
    │        │                ┌PRED            'be<[113:on]>[68:who]'       │
    │        │                │SUBJ         68┌PRED 'who']                  │
    │OBJ     │                │              ┌PRED    'on<[68:who], [144:Atripla]>'│
    │        │ADJUNCT {       │XCOMP-PRED    │SUBJ    [68:who]              │
    │        │                │           113│OBJ 144┌PRED 'Atripla']       │
    │        │                │PRON-REL      [68:who]                       │
    │        │              85│TOPIC-REL     [68:who]                       │
 11 └     44 └               └                                              ┘
```

# F-structures mapped to Abstract KR

*Find patients who were on Atripla.*



```
F-structure chart

"find patients who were on Atripla."

    PRED  'find<[11-SUBJ:null_pro], [44:patient]>'
    SUBJ  [PRED 'null_pro']
          [PRED      'patient'
          [                  [PRED          'be<[113:on]>[68:who]'
          [                  [SUBJ       68[PRED 'who']
          [                  [              [PRED     'on<[68:who], [144:Atripla]>'
    OBJ   [ADJUNCT {         [XCOMP-PRED    [SUBJ     [68:who]
          [                  [           113[OBJ  144[PRED 'Atripla']
          [                  [
          [                  [PRON-REL      [68:who]
    11[   44[             85[TOPIC-REL      [68:who]
```

AKR    subconcept(find-0, [find#v#1, detect#v#1, find#v#3])
       subconcept(Atripla-15, [drug_combo#n#1])
       alias(Atripla-15, [Atripla])
       subconcept(patient-3, [patient#n#1, affected_role#n#1])
       **role(ob, find-0, patient-3)**
       **role(prep(on), patient-3, Atripla-15)**

# AKR to Domain-Specific Logical Form

*Find patients who were on Atripla.*

**AKR**

subconcept(find-0, [find#v#1, detect#v#1, find#v#3]
subconcept(Atripla-15, [drug_combo#n#1]
alias(Atripla-15, [Atripla])
subconcept(patient-3, [patient#n#1, affected_role#n#1])
**role(ob, find-0, patient-3)**
**role(prep(on), patient-3, Atripla-15)**

**QUADRI**

**patient-has-drug-combo**(patient-3, Atripla-15)
sort(patient-3, patient)
sort(Atripla-15, drugCombo)

Plus quantifier information…

# Domain sort and relation mapping

- Domain relations have **argument signature**

  patient-has-regimen*(patient, regimen)*

  patient-has-test*(patient, medical_test)*

  regimen-has-drug-combo*(regimen, drug_combo)*

  test-time*(medical_test, time_point)*

  test-has-value*(medical_test, test_result)*

- Words (phrases) labeled for sort

  *patient_1:patient*

  *viral_load_2:medical_tes*t

# Task-specific customization: AKR → logical form

- Identification of quantifiers

    e.g. *every, some, at least 2, many*, …

- Cardinality vs. Measure specification

    e.g. *at least 2 regimens, at least 8 weeks*, …

- Mapping of conditions associated with quantified terms
    - e.g. distinguishing between restriction and nuclear scope
    - *every patient who has property P  (does X) :   ∀x ( patient(x) & P(x) → … )*
    - *every patient has property P :                    ∀x ( patient(x) → P(x) )*

- Fix scope relations between quantifiers
    - AKR underspecifies the scope of quantified terms
    - Scoping restrictions imposed by the grammar
    - Heuristics for fixing scope underdetermined by the grammar

# Task-specific customization: AKR → logical form

Donkey anaphora

*every patient on <u>a regimen </u>with AZT failed <u>the regimen </u>after 24 weeks*

Dependencies between terms do not align with syntactic structure

*a patient with norvir had a high viral load*

synlink(ob,have-6,viral_load_3)
synlink(sb,have-6,patient_8)

synlink(prep(with),patient_8,norvir_2)

synlink(nn_element,viral_load_3,high_7)

semlink(patient_8,regimen_9)
semlink(patient_8,viral_load_3)
semlink(regimen_9,norvir_2)
semlink(viral_load_3,high_7)

*every patient had a high viral load after 8 weeks on a regimen with norvir*

Interpretation in the domain

Disambiguation via interpretation in the domain

# Illustration

Find a **patient** on a **regimen** with **norvir and epivir who** had a **high viral load**

| | | | | | |
|---|---|---|---|---|---|
| Patient | Regimen | Drugs | Wh-word | Test-value | Medical test |

# Illustration



Find a **patient** on a **regimen** with **norvir and epivir who** had a **high viral load**

| Patient | Regimen | Drugs | Wh-word | Test-value | Medical test |

# Illustration



Find a **patient** on a **regimen** with **norvir and epivir who** had a **high viral load**

| Patient | Regimen | Drugs | Patient | Test-value | Medical test |

Patient-Has-Regimen

Regimen-has-drugs

Test-has-value

Patient-has-test

# Find linguistic links between word-pairs that match argument signatures

*Find patients who had a high viral load after 24 weeks on a regimen with Atripla.*

Direct Link (e.g. preposition)

> *role(prep(with), regimen_3, atripla_4)*
>
> →**semlink(regimen_3, atripla_4, via(prep(with))**

*record linguistic link*

*Coarguments of the verb*

> *role(sb, have, patient_1)*
>
> *role(ob, have, viral_load_2)*
>
> →**semlink(patient_1, viral_load_2, via(have))**

*record linguistic link*

# Semlinks map to DS relations

- Independent of linkage structure

  **semlink(X, Y, Z) iff  X:patient, Y:medical_test, Z:any**

  ➔  **patient-has-test(X, Y)**
  *a patient (had / with) a high viral load*

- Specific  to linkage structure

  **semlink(X, Y, Z)**
  **iff  X:time_period, Y:regimen,  Z:via(prep(on))**

  ➔**initial-interval(X,Y)**

  *24 weeks on a regimen  vs. 24 weeks after a regimen*

# Recovering implicit terms and relations

*patient-has-test(M, Test )*

➔

*exists time_point TP,*
*test-time(M, Test, TP)*
*-----------*
patient-has-drug(P, D)

➔

*exists regimen R,*
*patient-has-regimen(P, R)*
*regimen-has-drug(R, D)*

# Ambiguity management

Options multiplied out

$$\left\{ \begin{array}{l} \textit{The sheep-sg liked the fish-sg.} \\ \textit{The sheep-pl liked the fish-sg.} \\ \textit{The sheep-sg liked the fish-pl.} \\ \textit{The sheep-pl liked the fish-pl.} \end{array} \right\}$$

Options packed

$$\textit{The sheep} \left\{ \begin{array}{l} sg \\ pl \end{array} \right\} \textit{liked the fish} \left\{ \begin{array}{l} sg \\ pl \end{array} \right\}$$

Packed representation:
- – Encodes all dependencies without loss of information
- – Common items represented, computed once
- – Key to practical efficiency with broad-coverage grammars

# Packing

- Calculate and represent compactly all analyses at each stage
- Pass all or N-best analyses along through the stages
- Mark ambiguities in a free choice space
- Choice space:
- $\quad$ A1 $\vee$ A2 $\leftrightarrow$ *true*
- $\quad$ A1 $\wedge$ A2 $\rightarrow$ *false*

# Ambiguity passed on in AKR →LF mapping

*The patient  [had   [a regimen [with norvir]]]*

*[had    [a martini [with an olive]]]*

*The patient  [had  [a regimen]    [with norvir]]*

*[had   [a martini]   [with Olivia]]*

**Choice: (A1 xor A2) iff 1**

   **A1:** role(prep(with), have-1, norvir-5)

   **A2:** role(prep(with), regimen-12, norvir-5)

# Reducing choice space via selection

*(Regimen, Drug)*

**role(prep(with), R, D)**    *(R : regimen, D : drug)*

➡

*semlink(R, D,via(prep(with)))*

Semantically meaningful attachments eliminate uninterpretable choices

**role(prep(with), %%, D)**

➡

**stop.**

# Using interpretability to disambiguate

*The patient [had [a regimen [with norvir]]]*

~~*The patient [had [a regimen] [with norvir]].*~~

**Choice: (~~A1 xor~~ A2) iff 1**

    subconcept(have-1, [have#v#1, use#v#1, have#v#3]

    subconcept(norvir-5, [drug#n#1]

    alias(norvir-5, [norvir])

    subconcept(patient-3, [patient#n#1, affected_role#n#1])

    role(sb, have-1, patient-3)

    role(ob, have-1, regimen-12)

  ~~**A1:** role(prep(with), have-1, norvir-5)~~

  **A2: role(prep(with), regimen-12, norvir-5)**   *(Regimen, Drug)*

# Ambiguities multiply
# e.g. from prepositional attachment

***Find  patients who had a high viral load** after*
***24 weeks on a regimen with norvir.***

(62 ways ambiguous)

xor(A1, A2, A3, A4, A5, A6, A7, A8, A9, A10, A11, A12, A13, A14, A15, A16, A17, A18, A19, A20, A21, A22, A23, A24, A25, A26, A27, A28, A29, A30, A31, A32, A33, A34, A35, A36, A37, A38, A39, A40, A41, A42, A43, A44, A45, A46, A47, A48, A49, A50, A51, A52, A53, A54, A55, A56, A57, A58, A59, A60, A61, A62) iff 1

Conceptual Structure:
  or(A27,A29,or(or(A62,A61,A60,A59,A58,A57,A56,A55,A54,A53,A52,A51,A50,A49,A48),or(A47,A46,A45,A44,A43,A42,A41,A40),or(A35,A34),or(A33,or(A37,A36)),A39)):
    role(nn_element,viral_load-26,high-23,1)
    definite(regimen-47)
    subconcept(find-0,
      [find#v#1,detect#v#1,find#v#3,determine#v#1,find#v#5,witness#v#2,line_up#v#2,discover#v#2,discover#v#4,find#v#10,rule#v#4,receive#v#2,find#v#13,recover#v#1,find#v#15,find_oneself#v#1])
  or(A12,or(A14,or(or(A24,A23),A21),or(or(A26,A25),or(A28,A31))),or(A38,or(A45,A44),or(or(A55,A54),A52,or(or(A57,A56),A58,or(A62,A60)))))):
    role(prep(with),find-0,norvir-50)
  …

# Meaningful attachments

~~or(or(or(A28,A29),A24),A23):~~
~~role(prep(after),patient-7,24-22)~~
~~or(or(or(A42,A43),or(A38,A39),A33,A32),A10,A9):~~
~~role(prep(after),patient-7,week-26)~~
~~or(or(A30,A31),A22,A21):~~
~~role(prep(after),viral_load_test-16,24-22)~~

or(or(or(A41,or(A44,A45)),A40,or(A35,A36,A37),A34)…
    **role(prep(after),viral_load_test-16,week-26)**         *(Viral-load, Time-Period)*
        ➔ *medical-test-has-time(viral_load_test-16, time_point-1)*
        ➔ *occurs-after(time_point-1, week-26)*

or(A36,or(or(or(A27,A28,A29),…))
    **role(prep(during),week-26,regimen-37)**         *(Time-period, Treatment)*
        ➔ *occurs-during(week-26,regimen-37)*

~~or(or(A17,A18),or(or(A13,A14),A11,A12),A4,or(A1,A2,A3)):~~
~~role(prep(with),find-1,viral_load_test-16)~~

or(or(or(or(A28,A29),A24,….)
    **role(prep(with),patient-7,viral_load_test-16)**         *(Patient, Medical-Test)*
        ➔ *patient-has-test(patient-7, viral_load_test-16)*
                …)

# Meaningful analyses survive

(or(A32,or(or(or(or(A33,A34),A27,A25),A26….)….

    **role(cardinality_restriction,week-26,24)**            *(Time-Period, Cardinality)*

        ➔ *interval-has-duration(week-26,24 weeks)*

or(or(or(A45,or(A41,A42),A39)….)….

    **role(nn_element,viral_load_test-16,high-11,1)**       *(Medical-Test, Test-Value)*

        ➔ *medical-test-has-value(viral_load_test-16, high-11)*

or(or(or(A41,or(A44,A45)),A40,or(A35,A36,A37),A34)…

    **role(prep(after),viral_load_test-16,week-26)**       *(Viral-load, Time-Period)*

        ➔ *medical-test-has-time(viral_load_test-16, time_point-1)*

        ➔ *occurs-after(time_point-1, week-26)*

or(A36,or(or(or(A27,A28,A29),…)

    **role(prep(during),week-26,regimen-37)**       *(Time-period, Treatment)*

        ➔ *occurs-during(week-26,regimen-37)*

or(or(or(or(A28,A29),A24,….)

    **role(prep(with),patient-7,viral_load_test-16)**       *(Patient, Medical-Test)*

        ➔ *patient-has-test(patient-7, viral_load_test-16)*

# Bridge flattened LF given to reasoner

*Find patients who had a high viral load
after 24 weeks on a regimen with Atripla.*

*(desired_answer patient_3)*

*(exists patient_3 sort patient)*
*(exists regimen_4 sort regimen)*
*(scopes-over restriction patient_3  regimen_4)*
  *…*

*(in restriction patient_3 (patient-has-regimen patient_3 regimen_4)*

*(in restriction regimen_4 (regimen-has-drug-combo regimen_4 Atripla_2))*

*(in nscope patient_3 (patient-has-test-at-time patient_3 viral_load_5 high_1 time_point_7))*

*(in restriction time_point_7 (after time_point_7, week_6))*

*(in restriction week_6 (starts-at week_6 time_point_8))*

*(in restriction week_6 (starts-at regimen_4 time_point_8))*

*(in restriction week_6 (time_measure week_6 24 week))*

# Complex queries: multiple questions

- *Find patients who had a high viral load after 24 weeks on a regimen with Atripla;*

- *the patients exhibited M184v near the end of the regimen;*

- *the patients switched to a salvage regimen with boosted EFV.*

# Points to remember

- Experts use many abstractions over information in DB

- A reasoner can link higher level abstractions found in natural queries with combinations of data base elements

- Mapping language in a specific domain can be guided by signatures of higher level domain relations
(only sometimes requiring specific constructs)

- Mapping to domain relations can be used to eliminate uninterpretable ambiguities

# Porting to a new domain

- Requires being able to build a language model for that domain.

- This was tried in the intelligence community (IC) domain

- RDF class definition triples were used as our *argument signatures* over the existing Quadri system.

# Terrorism

- 



Query: Which member of the Egyptian Islamic Jihad killed more than 30 people in Cairo and injured the president's daughter?

Query Term and Value    RDF Relation    Domain Rewrite

**RDF INFO**    Arguments    Matching pair in NL Question    **Relation specification**
(defquery-template !mrrt:**HumanAgentKillingAPerson**    (evenLocationGPE'
 (!mrqv:**Event** !mric:killingHumanAgent !mrqv:**killingHumanAgen**t) **(kill, member, via(sb))**     (Event, GPE, restrict(none)))
 (!mrqv:**Even**t !rdf:type !mric:**HumanAgentKillingAPerson**)) ...    (PersonHasDaughter'
     (Person, daughter, restrict(poss))
(defquery-template !mrrt:**eventLocationGPE**    Type/Concept Info    ....
 (!mrqv:**Event** !mric:eventLocationGPE !mrqv:**GPE**)) ...    **(kill, Cairo, via(prep(in)))**

    Link info    **Named Concepts**
(defquery-template !mrrt:**HumanOrganizationHasMember**    (person'
 (!mrqv:**HumanOrganization** !mric:hasMember !mrqv:**member**) **(egyptian_jihad, member, via(prep(of)))**     (daughter, member, agent,...))
 (!mrqv:**member** !rdf:type !mric:**Person**)) ...    (daughter'
     (daughter))
(defquery-template !mrrt:**PersonHasDaughte**r    (event'
 (!mrqv:**Perso**n !mric:hasDaughter !mrqv:**daughter**)    **(president,daughter, via(poss))**     (killingEvent,InjureEvent,...))
 (!mrqv:**daughte**r !rdf:type !mric:**Person**))    (killingEvent'
     (HumanAgentKillingAPerson))
(defquery-template !mrrt:**HumanInjuryEventPersonInjured**    ....
 (!mrqv:**Even**t !mric:personInjured !mrqv:**Person**)    **Partly integrated into BRIDGE**
 (!mrqv:**Event** !rdf:type !mric:**HumanInjuryEven**t))    **(injure,daughter, via(ob))**     **ontology**

# System Output

((input In July a terrorist from Algeria who is associated with Al Qaeda killed nearly 30 people in Yemen.)
(quant exists Algeria_1 sort geo_political_entity)
(quant exists Yemen_2 sort geo_political_entity)
(quant exists terrorist_3 sort person)
(quant exists Al_Qaeda_4 sort human_organization)
(quant (complex_card nearly 30) people_5 sort people_group)
(quant exists July_6 sort date)
(definite Algeria_1)
(definite Yemen_2)
(**HumanAgentKillingAPerson** (kill:n 164 terrorist_3)
(**HumanKillingEventPersonKilledUpperBound** (kill:n 164 people_5)
(**eventLocationGPE** (kill:n 164 Yemen_2)
(**organizationHasMember** Al_Qaeda_4 terrorist_3)
(**PersonHasBirthPlace** terrorist_3 Algeria_1)
(**EventHasDate** (kill:n 164 July_6)
)

# Thank you.