# Towards Declarative and Efficient Querying on Protein Structures
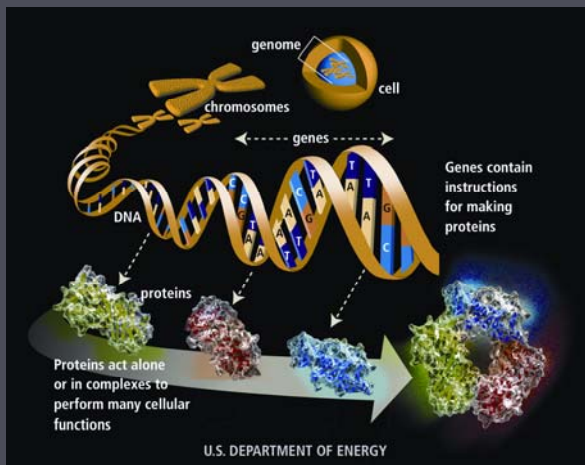
Jignesh M. Patel

University of Michigan

---

# Biology 101



When there is something to do in the cell, it is a protein that *does* it.

## Data Types

Sequences: AGCGGTA….

Structure:

Interaction Maps:

Micro-arrays:

Gene A
Gene B

# Role of DBMS in Bioinformatics

► Large Data Sets
  ▪ Growing exponentially!
► Data Types
  ▪ Sequences/arrays/3-D/text
► Complex Queries
  ▪ Ad-hoc tools today
  ▪ Integrated querying done using procedural methods
► Scalability/Parallelism
  ▪ Home-grown techniques often used today
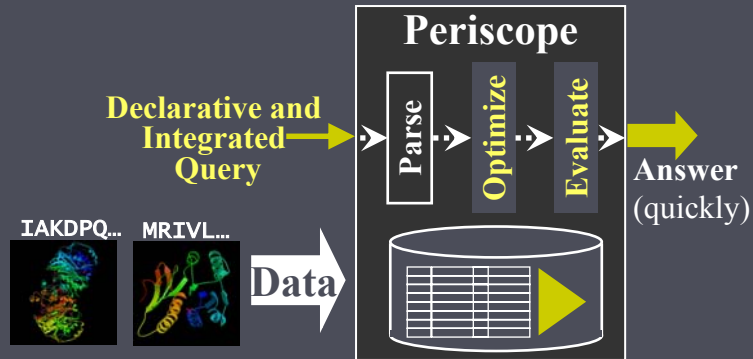► Need: Declarative and efficient querying



Source: GenBank

# Protein Querying

► Protein functionality determined by
  ▪ Sequence composition
  ▪ Geometric structure
► Need to query across all structures
► Current methods
  ▪ Non-declarative tools, often not very efficient on large data sets
  ▪ Support querying on only one structure

# Periscope

**Goal: Design, implement, and evaluate a database management system for <u>declarative</u> and <u>efficient</u> querying on <u>all</u> protein structures**

# Roadmap

- ► Background and Introduction
- ► Primary Structure Sequence Matching
- ► Querying on Secondary Structure
- ► PiQA: Integrated Query Algebra
- ► Summary

# Sequence Matching

▶ Find similar sequences
  - Given a protein sequence, find homologous matches in the database
  - Similarity based on "local similarity"
  - A local-alignment algorithm
  - Operations: Replace, Delete, Insert
  - Score using a substitution matrix

**<u>Database:</u>** THE T RAIN DRIVER'S CABIN

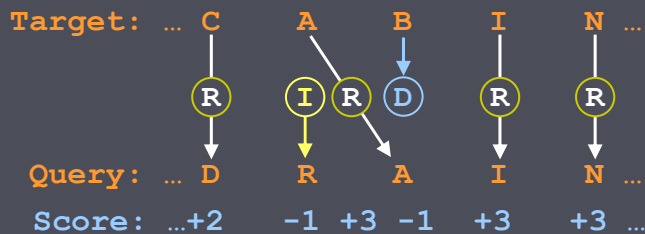**<u>Query:</u>** D RAIN

---

# Sequence Matching

▶ Find similar sequences
  - Given a protein sequence, find homologous matches in the database
  - Similarity based on "local similarity"
  - A local-alignment algorithm
  - Operations: Replace, Delete, Insert
  - Score using a substitution matrix

Target

|   | -  | A  | B  | C  | D  |
|---|----|----|----|----|----|
| **-** | -1 | -1 | -1 | -1 | -1 |
| **A** | -1 | 3  | 0  | 0  | 0  |
| **B** | -1 | 0  | 3  | 0  | 0  |
| **C** | -1 | 0  | 0  | 3  | 0  |
| **D** | -1 | 0  | 0  | 2  | 3  |

Query

```
Target:  … C      A      B      I      N …
            (R)   (I)(R) (D)    (R)    (R)
Query:   … D      R      A      I      N …
Score:   …+2     -1 +3 -1      +3     +3 …
```

# Smith-Waterman

**C A B I N**

S-W Matrix: G

|   | C | A | B | I | N |
|---|---|---|---|---|---|
| **D** | 2→1 |  |  |  |  |
| **R** | 1 | 2→1 |  |  |  |
| **A** |  | 4→3→2→1 |  |  |  |
| **I** |  | 3 | 4 | 6→5 |  |
| **N** |  | 2 | 3 | 5 | 9 |

**Target**

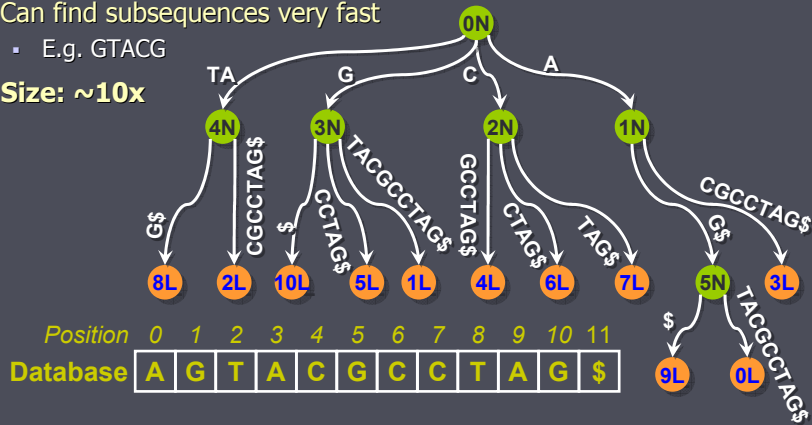|   |   | - | A | B | C | D |
|---|---|---|---|---|---|---|
| Query | - | -1 | -1 | -1 | -1 | -1 |
|  | A | -1 | 3 | 0 | 0 | 0 |
|  | B | -1 | 0 | 3 | 0 | 0 |
|  | C | -1 | 0 | 0 | 3 | 0 |
|  | D | -1 | 0 | 0 | 2 | 3 |

Substitution Matrix

$$G_{i,j} = \max \begin{cases} 0 \\ G_{i-1,\,j-1} + S(q_i \rightarrow t_j) \\ G_{i-1,\,j} + S(q_i \rightarrow -) \\ G_{i,\,j-1} + S(- \rightarrow t_j) \end{cases}$$

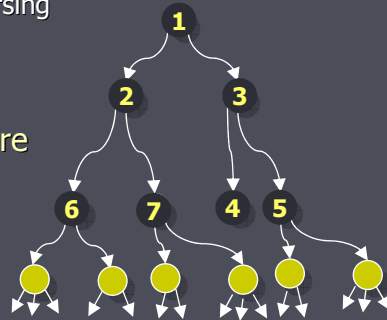C → R → D
A → I R → R
B → D → A
I → R → I
N → R → N

---

# Suffix Trees

- Compact Patricia trie
  - Every suffix has a path (to leaf)
  - Every subsequence is a prefix of a path
- Can find subsequences very fast
  - E.g. GTACG
- **Size: ~10x**
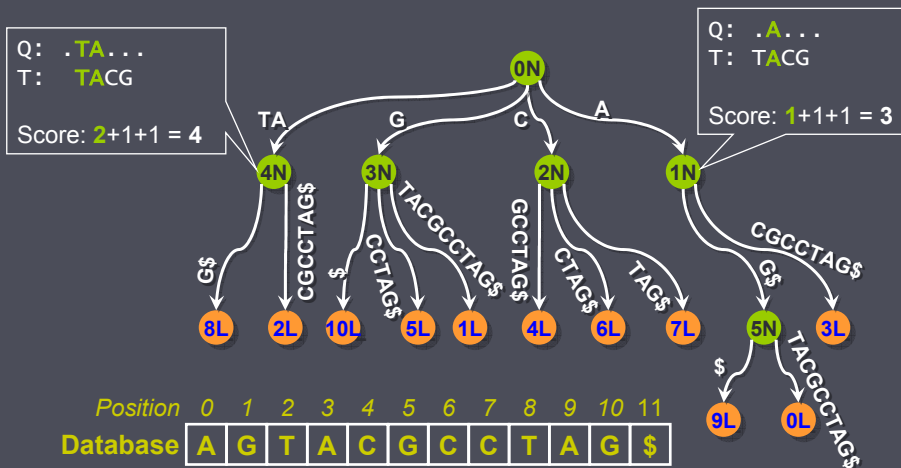
| Position | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Database** | A | G | T | A | C | G | C | C | T | A | G | $ |

# OASIS

- ► Search driven by the suffix tree
  - ▪ Fill up the S-W columns by traversing down the suffix tree
- ► Best-first: expand node in the tree with highest expected score
  - ▪ Expected Score = current score + best possible score for unconsumed portion of query
  - ▪ Guarantees **online** behavior!
- ► Exploits redundancy in the database

2/7/2004 © Jignesh M. Patel, 2004 11



# OASIS

- ► Query: TACG
- ► Unit Edit Distance Matrix: Same Symbol Substitution = 1, else -1

Q: .TA...
T: TACG

Score: 2+1+1 = 4

Q: .A...
T: TACG

Score: 1+1+1 = 3

Position  0  1  2  3  4  5  6  7  8  9  10  11
Database  A  G  T  A  C  G  C  C  T  A  G  $

# Experimental Setup: Comparison

- ► Data set: swissprot
  - ▪ 110K proteins, # symbols: 40M, data size: 40MB
  - ▪ Index Size: 500MB (~12.5 bytes per symbol)
- ► Workload: 100 queries from ProClass motif database
  - ▪ Short queries – lengths 6 - 56, avg len = 16
- ► PAM30 scoring matrix
- ► BLAST parameters  (short query settings, 5-15 residues)
  - ▪ E=20,000, word size = 2
- ► OASIS: $minScore = \left\lceil \dfrac{\ln(Kmn) - \ln(E)}{\lambda} \right\rceil$

- ► Platform: 1.7GHz Xeon, Linux, 256MB buffer pool, 2K page size, clock replacement policy

---

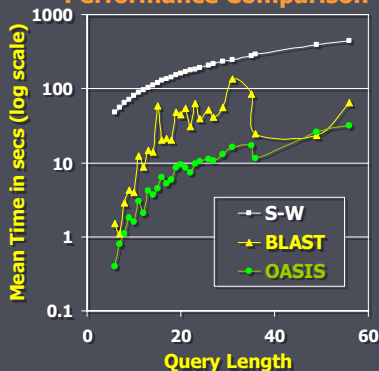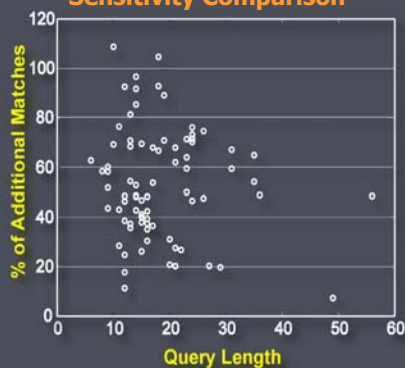# Experimental Results

**Performance Comparison**



**Sensitivity Comparison**



- ► OASIS 1-2 orders of magnitude faster than S-W, and usually also faster than BLAST
- ► OASIS retrieved ~ 60% more matches than BLAST

**Lots more in our VLDB'03 paper**

# Roadmap

- ▶ Background and Introduction
- ▶ Primary Structure Sequence Matching
- ▶ Querying on Secondary Structure
- ▶ PiQA: Integrated Query Algebra
- ▶ Summary

---

# Querying Protein Secondary Structure

- ▶ Secondary Structure describes protein folds
  - Beta sheets (e); Alpha helices (h); Loops (l)
  - Predicted structure; helps determine protein function
- ▶ Data Model
  - Sequence of Segments
  - 'h h h l l e e e e' → <3h>, <2 l>, <4 e>
- ▶ Query Language
  - Sequence of regular expression terms
    - ▶ Example: Beta sheet of length 4-10 followed at some point by an helix of length 3-6
    - **Query: <e 4 10> <? 0 Inf> <h 3 6>**

# Schema and Query Evaluation

| id | name | len | sec-seq | ... |
|----|------|-----|---------|-----|
| 1  | A    | 5   | lleee   | ... |
| 2  | B    | 6   | hhhee   | ... |

**Protein Table**

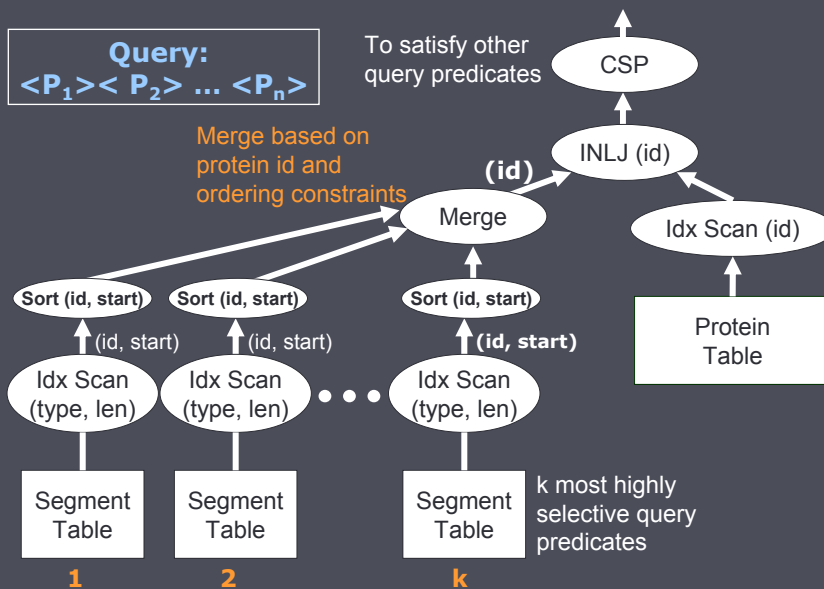| seg-id | id | type | len | start-pos |
|--------|----|------|-----|-----------|
| 1      | 1  | l    | 2   | 1         |
| 2      | 1  | e    | 3   | 3         |
| 3      | 2  | h    | 3   | 1         |
| 4      | 2  | e    | 2   | 4         |

**(Derived) Segment Table**

## Evaluation Methods

- ► CSP: complex scan
- ► SSS: scan segment table
- ► ISS: use segment index
- ► MISS (k): multiple ISS

Scan each protein tuple and evaluate query using a finite state automata

N-way merge join of **k** multiple segment index scans + FK join

---

# Multiple Index Scan Method: MISS(k)



**Query:**
**<P₁><P₂> ... <Pₙ>**

To satisfy other query predicates — CSP

Merge based on protein id and ordering constraints

**(id)** — INLJ (id)

Merge

Idx Scan (id)

Sort (id, start) — Sort (id, start) — Sort (id, start)

(id, start) — (id, start) — **(id, start)**

Idx Scan (type, len) — Idx Scan (type, len) • • • Idx Scan (type, len)

Protein Table

Segment Table — Segment Table — Segment Table

k most highly selective query predicates

**1**     **2**     **k**

# Query Optimizer

- ► Chooses best method for given query
- ► Optimization:
  - Requires estimates of query predicate selectivities and result cardinality
  - Cost model: CPU and IO
- ► Estimation: Two types of histograms
  - Basic: segment predicate selectivity
  - Complex: query selectivity

---

# Basic Histogram

- ► Estimates selectivity of individual query predicates
- ► 2-D array
  - Dimensions : Length, Fold type
  - Value: number of each <type,len> segment
- ► Example
  - Pred <e 4 4>, Est: 52
  - Pred <e 2 4>, Est: 35+45+52

| Len | H | E | L |
|-----|-----|-----|-----|
| 1 | 20 | 10 | 18 |
| 2 | 23 | 35 | 25 |
| 3 | 36 | 45 | 33 |
| 4 | 44 | 52 | 35 |
| ... | | | |
| 50 | 4 | 2 | 0 |
| - | 1 | 0 | 0 |

# Complex Histogram

- ▶ Estimates **result** cardinality
- ▶ Four-dimensional structure
  - Protein id (equi-width buckets)
  - Start position (equi-width buckets)
  - Length (1 – 50) ⎤ Same as in the
  - Type ('e', 'h', 'l') ⎦ basic histogram
- ▶ Example:  Position [x] [y] [z] ['e']
  - holds the number of <e z z> segments whose starting position is in the range of the $y^{th}$ bucket and whose protein id lies within the $x^{th}$ bucket range
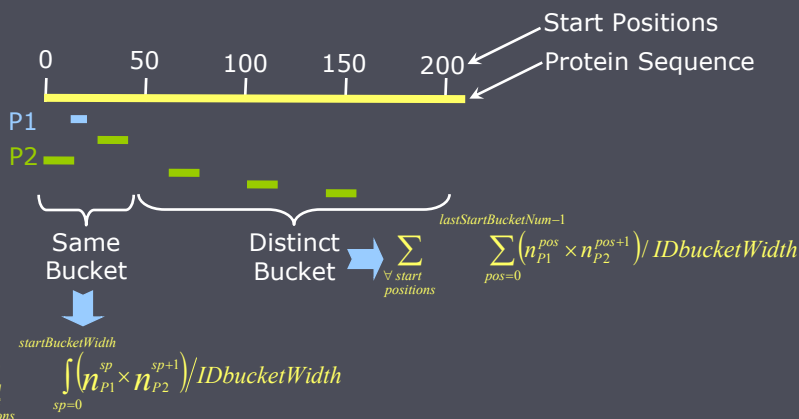
---

# Complex Histogram

4 Dimensions
- Protein id
- Start pos
- Length
- Type

- ▶ Query: {<P1> <P2>}
  - Compute selectivity of query result

Start Positions

0    50    100    150    200 ← Protein Sequence

P1

P2

Same Bucket

Distinct Bucket → $\sum_{\forall\ start\ positions} \sum_{pos=0}^{lastStartBucketNum-1} \left( n_{P1}^{pos} \times n_{P2}^{pos+1} \right) / IDbucketWidth$

$\sum_{\forall\ start\ positions} \int_{sp=0}^{startBucketWidth} \left( n_{P1}^{sp} \times n_{P2}^{sp+1} \right) / IDbucketWidth$
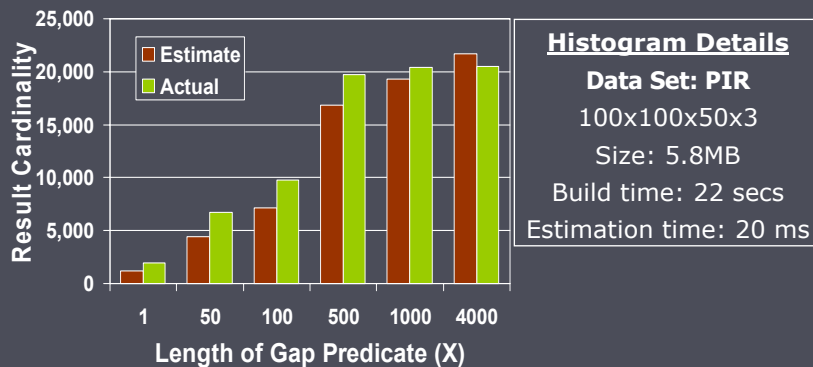
# Experimental Evaluation

- ► Techniques implemented in Periscope
  - ▪ Configuration: Using the SHORE storage manager, 64MB buffer pool size, 16K page size
- ► Also implemented in a commercial ORDBMS
  - ▪ Used type extensibility to create array-like data types for sequences & user-defined functions for FSM
- ► Data Set: PIR
  - ▪ 250K protein tuples, 250MB
  - ▪ Segment Table: 10M tuples, 355MB
- ► Platform: 1.7GHz Xeon, Linux, 40GB SCSI disk

---

# Complex Histogram Accuracy

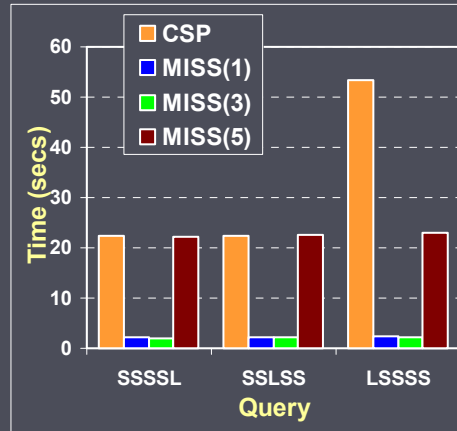- ► Query: {<l 15 15> <? 0 X> <h 24 24>}
- ► Experiment: Vary Gap Predicate



**Histogram Details**
**Data Set: PIR**
100x100x50x3
Size: 5.8MB
Build time: 22 secs
Estimation time: 20 ms

Histogram accurate within 80% of the actual result

# Query Performance

- ▶ Query: 9 Predicates
  $P_1$ $G_1$ $P_2$ $G_2$ $P_3$ $G_3$ $P_4$ $G_4$ $P_5$
  P=predicate, G=Gap pred.
- ▶ Expt: Vary Selectivity of P's
  Alternatives: S (0.3%), L (7%)

- ▶ Choice of algorithm is critical
- ▶ CSP sensitive to position of **L** pred.
- ▶ Choice of k in MISS critical
  Index Probe
  + reduce # protein tuples fetched
  - probe cost, sorting and merging
  Influenced by query and predicate
  selectivity
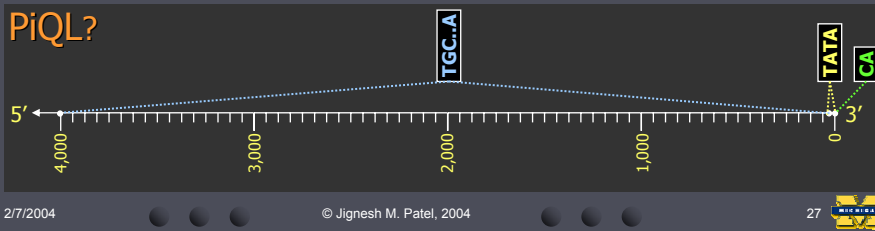


**More details in our VLDB'02 paper**

---

# Roadmap

- ▶ Background and Introduction
- ▶ Primary Structure Sequence Matching
- ▶ Querying on Secondary Structure
- ▶ PiQA: Integrated Query Algebra
- ▶ Summary

## Slide 1

# PiQA

- ► An algebra in PNF for queries on primary & secondary structures
- ► Examples:
  - Match the primary sequence "AAANBPPPPSDF", but ignore mismatch in the segment "NBPPP" if it is on a loop
    (P.p * "AAA") || ((P.p * "NBPPP") U (P.s * <L 5 5>)) || (P.p * "PSDF")
  - BLAST: Match n-grams + match extension + transitive closure
  - Find all occurrences of 'TGCTGACTCAGCA' within 4000 bps upstream of a 'CA' with 'TATA' 25-30 bps upstream of the 'CA'
    M * 'TGCTGACTCAGCA' || $_{3957}$ M * 'TATA' ||$_{26}$ M * 'CA'

PiQL?

TGC.A

TATA  CA

5'

3'

4,000    3,000    2,000    1,000    0

---

## Slide 2

# Roadmap

- ► Background and Introduction
- ► Primary Structure Sequence Matching
- ► Querying on Secondary Structure
- ► PiQA: Integrated Query Algebra
- ► Summary

---

# Summary

- ► Bioinformatics applications (urgently) need <u>declarative</u> and <u>efficient</u> query processing tools
  - ▪ Current procedural methods reminiscent of pre-relational days
- ► Database researchers have a lot to contribute!
- ► Periscope is our effort in this direction
  - ▪ PiQA: algebraic framework for querying on primary and secondary structures
  - ▪ Primary Structure: OASIS
  - ▪ Declarative querying on secondary structure: Periscope/PS$^2$
- ► Current Status
  - ▪ OASIS and Periscope/PS$^2$ currently (beta-)deployed at UM
  - ▪ Under development: "PiQA powered PiQL queries on Periscope"
- ► http://www.eecs.umich.edu/periscope