# $D_o$MiNO - Spatial data mining exploring co-location of adverse birth outcomes and environmental variables.

**Dr. Alvaro R. Osornio-Vargas**
**Canadian Perinatal Programs Coalition**
**June 23rd & 24th, 2014**
**Montréal, Quebec**

# D$_o$MiNO is:

- Interdisciplinary and Exploratory
- Uses publicly funded data bases
- DATA MINING

# Background

Current research identifies *associations between ABO and various determinants of health:*

- – social factors (e.g. poverty, stress),
- – biological factors (e.g. diabetes, infection, maternal age)
- – environmental pollutants (e.g. metals, $PM_{10}$, $SO_2$)

# Growing evidence linking Urban Air Pollutants and Adverse Birth outcomes

# Background

Complex problem:

- Multiple sources of pollutants (e.g. traffic, industry).

- Interactions, dispersion, transport and fate of pollutants

- Intrinsic toxicity of pollutants

- Interactions between, social, biological, chemical and physical factors

# Background

In order to advance research on links between _environmental pollutants and ABO,_ methods for comprehensive assessment of the multiple variables interacting in complex ways are required.
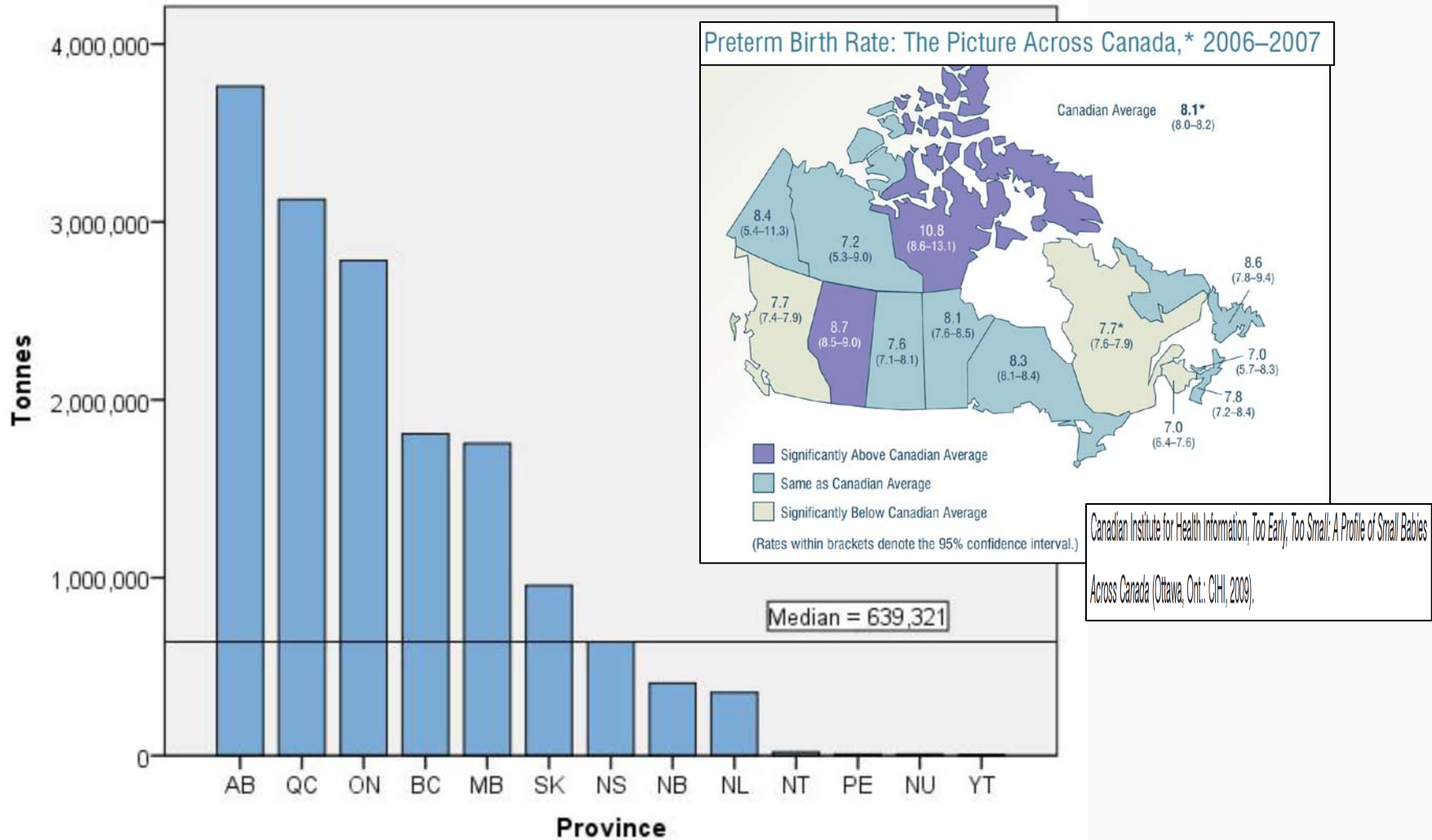
## DATA MINING

Spatial data mining exploring co-location of adverse birth outcomes and environmental variables. Osornio-Vargas, A.R.



Preterm Birth Rate: The Picture Across Canada,* 2006–2007

Canadian Average 8.1*
(8.0–8.2)

8.4 (5.4–11.3)
7.2 (5.3–9.0)
10.8 (8.6–13.1)
8.6 (7.8–9.4)
7.7 (7.4–7.9)
8.7 (8.5–9.0)
7.6 (7.1–8.1)
8.1 (7.6–8.5)
7.7* (7.6–7.9)
8.3 (8.1–8.4)
7.0 (5.7–8.3)
7.8 (7.2–8.4)
7.0 (6.4–7.6)

Significantly Above Canadian Average
Same as Canadian Average
Significantly Below Canadian Average

(Rates within brackets denote the 95% confidence interval.)

Median = 639,321

Canadian Institute for Health Information, Too Early, Too Small: A Profile of Small Babies Across Canada (Ottawa, Ont.: CIHI, 2009).
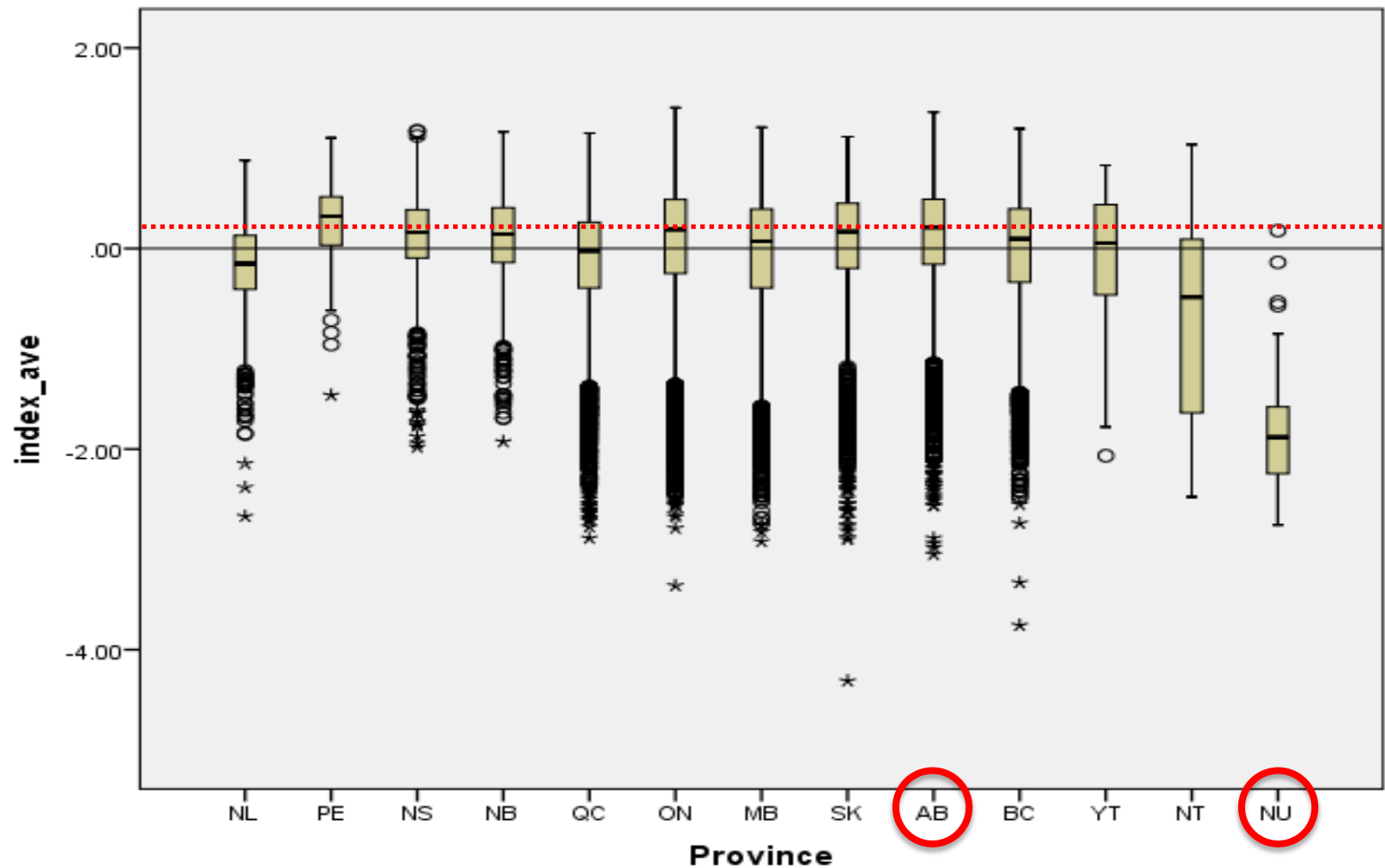
**Figure 2.** Tonnes of developmental toxicants reported to NPRI as released to air in each Province of Canada, from 2006 – 2010. Numbers of chemicals released vary by Province. The line marks median emissions.

# Distribution of SES index by province and territory

# Adverse birth outcomes and the environment
## 2002 - 2010

**Statistics Canada**
**National Pollutants Release Inventory**
**Wind Patterns**
**The Canadian Neonatal Network**
**Alberta Perinatal Health Program**

↓

**Data Mining**

↓

**Identify patterns**

↓

**Hypothesis**
**(e.g. collocation)**

# Objective

Aims to:

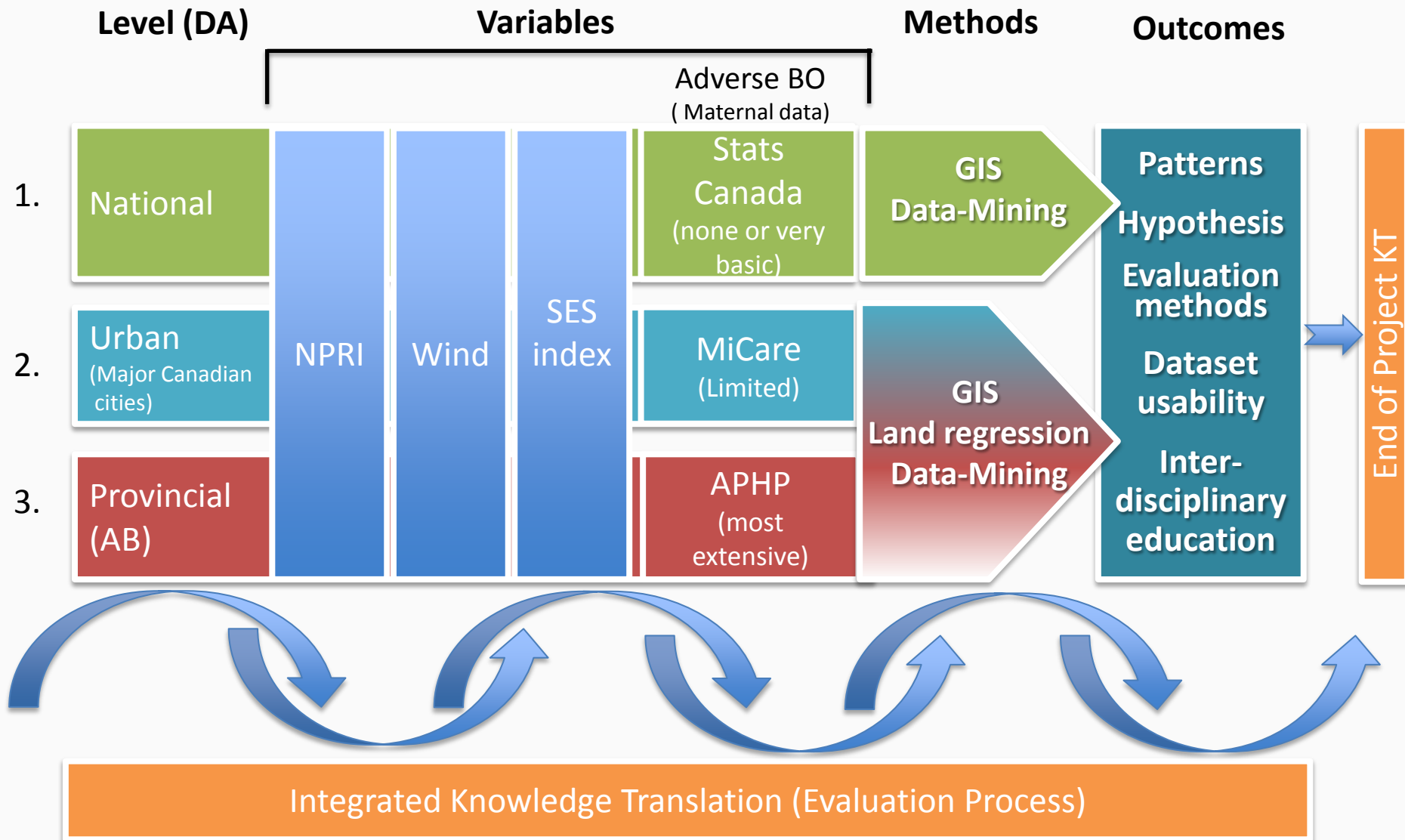Identify patterns

Generate <u>hypothesis</u>

Spatial collocation

Adverse birth outcomes

(pre-term, low birth weight and perinatal mortality)

Environmental Factors

(pollution, SES)

# Research Framework

| Level (DA) | Variables | | | Adverse BO ( Maternal data) | Methods | Outcomes | |
|---|---|---|---|---|---|---|---|
| 1. National | NPRI | Wind | SES index | Stats Canada (none or very basic) | GIS Data-Mining | Patterns Hypothesis Evaluation methods Dataset usability Inter-disciplinary education | End of Project KT |
| 2. Urban (Major Canadian cities) | | | | MiCare (Limited) | GIS Land regression Data-Mining | | |
| 3. Provincial (AB) | | | | APHP (most extensive) | | | |

Integrated Knowledge Translation (Evaluation Process)

*DoMiNO*
*Data Mining newborn outcomes*

| | Chemical name | Tonnes | | Chemical name | Tonnes |
|---|---|---|---|---|---|
| 1 | Sulphur dioxide | 7,614,400 | 30 | HCFC-22 | 487 |
| 2 | Carbon monoxide | 4,744,224 | 31 | Chloroform | 475 |
| 3 | Volatile Organic Compounds (VOCs) | 1,279,186 | 32 | Naphthalene | 449 |
| 4 | PM Total Particulate Matter | 962,176 | 33 | Arsenic (and its compounds) | 353 |
| 5 | $PM_{10}$ Particulate Matter $\leq$ 10 $\mu$m | 520,352 | 34 | Methyl methacrylate | 302 |
| 6 | $PM_{2.5}$ Particulate Matter $\leq$ 2.5 $\mu$m | 277,572 | 35 | Acetonitrile | 197 |
| 7 | Methanol | 69,679 | 36 | Tetrachloroethylene | 196 |
| 8 | n-Hexane | 26,108 | 37 | 1,3-Butadiene | 193 |
| 9 | Xylene (all isomers) | 25,897 | 38 | tert-Butyl alcohol | 126 |
| 10 | Toluene | 21,220 | 39 | Cadmium (and its compounds) | 121 |
| 11 | Hydrogen fluoride | 16,984 | 40 | Acrylonitrile | 69 |
| 12 | Carbon disulphide | 16,377 | 41 | Butyl benzyl phthalate | 47 |
| 13 | Styrene | 9,522 | 42 | N,N-Dimethylformamide | 46 |
| 14 | Methyl ethyl ketone | 8,653 | 43 | Sodium nitrite | 45 |
| 15 | Isopropyl alcohol | 6,947 | 44 | Benzo(a)pyrene - PAH | 42 |
| 16 | Acetaldehyde | 5,117 | 45 | Bis(2-ethylhexyl) phthalate | 37 |
| 17 | Ethylbenzene | 4,055 | 46 | 1,2,4-Trichlorobenzene | 27 |
| 18 | Benzene | 3,257 | 47 | p-Dichlorobenzene | 27 |
| 19 | Phenol (and its salts) | 3,031 | 48 | Ethylene oxide | 22 |
| 20 | 2-Butoxyethanol | 2,747 | 49 | Mercury (and its compounds) | 22 |
| 21 | Chloromethane | 2,242 | 50 | Biphenyl | 19 |
| 22 | Chlorine dioxide | 2,118 | 51 | Vinyl chloride | 10 |
| 23 | Methyl isobutyl ketone | 1,412 | 52 | Dibutyl phthalate | 7 |
| 24 | Trichloroethylene | 1,270 | 53 | 1,2-Dichloroethane | 6 |
| 25 | Lead (and its compounds) | 1,144 | 54 | 2-Ethoxyethyl acetate | 6 |
| 26 | Nickel (and its compounds) | 1,131 | 55 | Ethylene thiourea | 4 |
| 27 | Ethylene glycol | 830 | 56 | Bromomethane | 1 |
| 28 | Acrolein | 715 | 57 | Chlorobenzene | 1 |
| 29 | N-Methyl-2-pyrrolidone | 673 | 58 | Ethyl acrylate | 1 |
| | **Total** | **15,629,039** | | **Total** | **3,338** |
| | | | | **GRAND TOTAL** | **15,632,377** |

**Table I : Total amounts of developmental toxicants reported to NPRI as released to air in Canada in 2006 – 2010**
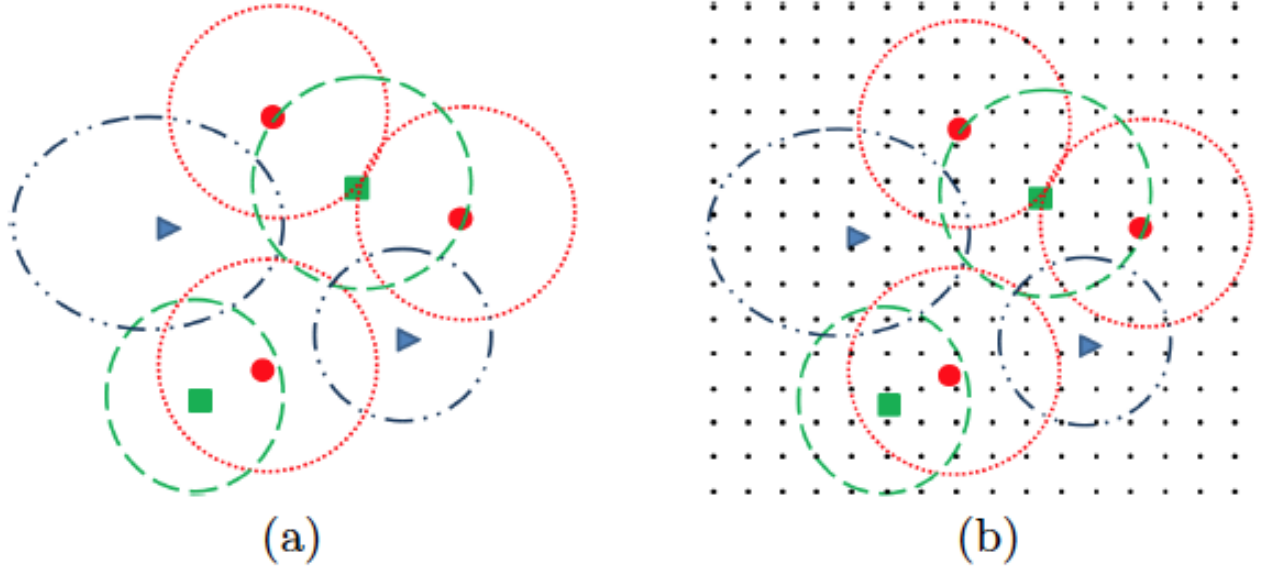
Fig. 1. Transactionization step: (a) An example spatial dataset with point feature instances and their buffers; (b) Grids imposed over the space.

$$p = \sum_{i=\sigma(XA)}^{\sigma(A)} \binom{n}{i} (P(X)P(A))^i (1 - P(X)P(A))^{n-i} \qquad (1)$$

**Algorithm 1** CMCStatApriori Algorithm.

**Require:** Set of antecedent features $F\backslash A$, the consequent feature $A$, derived transaction dataset $T$, the threshold $z_{min}$ for the $z$-score
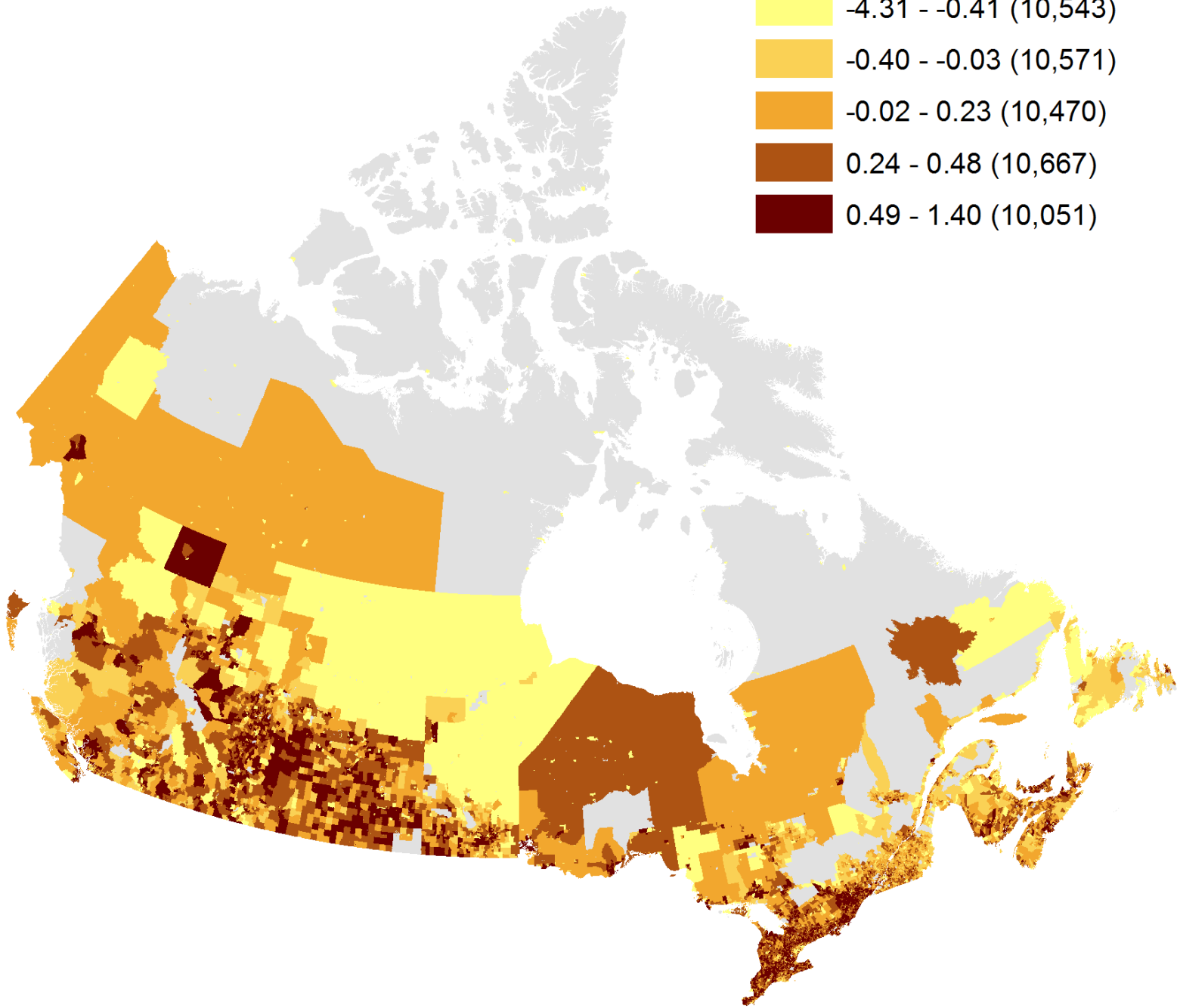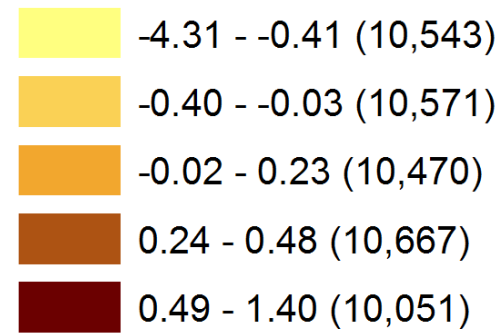
**Ensure:** Set of potential statistical significant co-location rules $P$

1: $P_1 = \{f_i \in F\backslash A | PS(f_i) = 1\}$
2: $l = 1$
3: **while** $(P_l \neq \varnothing)$ **do**
4: $\quad C_{l+1} = GenCands(P_l, A)$
5: $\quad P_{l+1} = PrunCands(C_{l+1}, z_{min}, A)$
6: $\quad l = l + 1$
7: **end while**
8: $P = \cup_i P_i$
9: **return** $P$

**SESIndex at the DA level**

- -4.31 - -0.41 (10,543)
- -0.40 - -0.03 (10,571)
- -0.02 - 0.23 (10,470)
- 0.24 - 0.48 (10,667)
- 0.49 - 1.40 (10,051)

Edmonton (ppb)

High : 28

Low : 0

▲ NAPS Monitors

0    1.5    3         6         9
Kilometers

# Table II. Evaluation of the integrated KT plan.

| Evaluation question | Proposed methods |
|---|---|
| *What perspectives do researchers from different disciplines and knowledge users have on interdisciplinary partnership research? How do these perspectives change over the course of the research project?* | 1. Participant observation<br>2. Individual semi-structured, time series interviews |
| *What challenges and barriers are experienced to interdisciplinary team development?* | 3. Participant observation<br>4. Individual semi-structured, time series interviews<br>5. Content analysis<br>6. Analysis of project log<br>7. End of project survey |
| *What strategies are most useful in building collaboration and addressing identified barriers?* | 8. Participant observation<br>9. Individual semi-structured, time series interviews<br>10. End of project survey |
| *What are the added benefits of interdisciplinary research-knowledge user collaboration?* | 11. End of project survey<br>12. End of project focus group<br>13. Analysis of project deliverables |

# Team

## University of Alberta

### Faculty of Medicine & Dentistry
| | |
|---|---|
| Dr. Osornio-Vargas* | Principal Investigator |
| Dr. Irena Buka | Children's environmental health |
| Dr. Khalid Aziz | Neonatology |
| Dr. Manoj Kumar | Neonatology |
| Dr. Sue Chandra | Obstetrics & Gynecology |
| Osnat Wine | Knowledge Translation |
| Emily Chan | Socioeconomic variables |

### Computing Sciences
| | |
|---|---|
| Dr. Osmar Zaiane* | Principal Investigator |
| Jundong Li | Data mining |
| Dr. Dr. Sajib Barua | Data mining |

### School of Public Health
| | |
|---|---|
| Dr. Sarah Bowen | Knowledge Translation |
| Dr. Yan Yuan | Biostatistics |
| Dr. Yutaka Yasui | Biostatistics |
| Jesus Serrano | Geostatistics |

### Faculty of Sciences
| | |
|---|---|
| Charlene Nielsen | Geostatistics |

## Carlton University

### Department of Health Sciences,
| | |
|---|---|
| Dr. Paul Villeneuve | Epidemiology |

## University of Victoria

### Interdisciplinary Studies
| | |
|---|---|
| Dr. Laura Arbour | Paediatrics and Genetics |
| Anders Erickson | GIS |

## Oregon State University

### School of Biological & Population Health Sciences
| | |
|---|---|
| Dr. Perry Hystad | Spatial exposure assessment |

## CAREX
| | |
|---|---|
| Dr. Eleanor Setton | Exposure Assessment |
| Dr. Paul Demers | Epidemiology |

## CIHR Maternal-Infant Care (MiCare) Program
| | |
|---|---|
| Dr. Prakeshkumar Shah | Neonatology |

## Knowledge Users

### Health Canada
| | |
|---|---|
| Dr. David Stieb | Epidemiology |
| Dr. Phil Blagden | Science Advisor |

### Alberta Perinatal Health Program
Nancy Aelicks

### Canadian Partnership for Children's Health & Environment
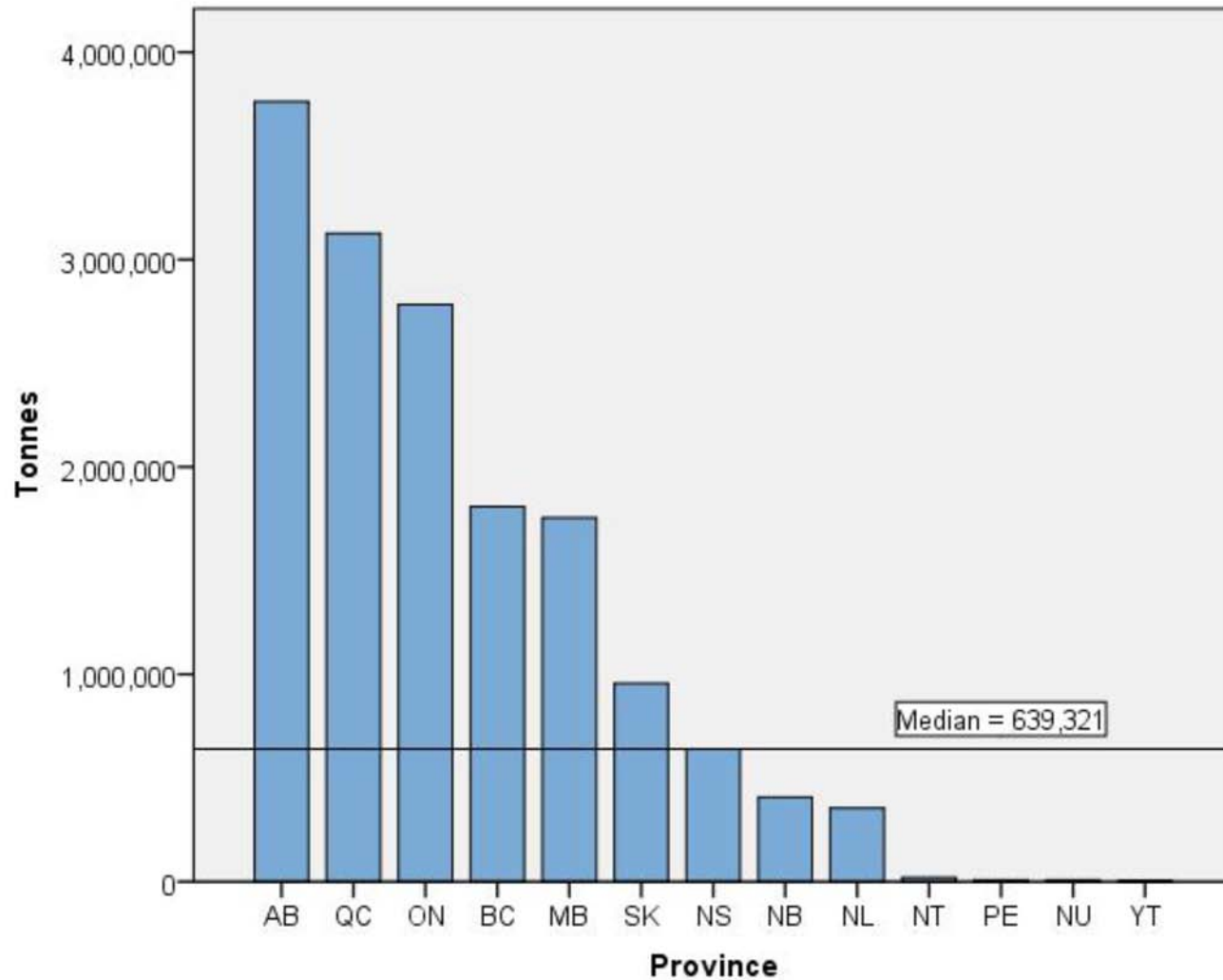Erica Phipps

# Thank you

osornio@ualberta.ca

**Figure 2.** *Tonnes of developmental toxicants reported to NPRI as released to air in each Province of Canada, from 2006 – 2010. Numbers of chemicals released vary by Province. The line marks median emissions.*

# Data Mining

## Conclusion

- New framework which uses (buffer & grid)-based transactionization and preserves spatial information better

- Statistical testing eliminates the usage of one global prevalence threshold

- Consideration of wind data and pollutants amounts to improve accuracy of results

25