# Visualization for Statistical Term Network in Newspaper
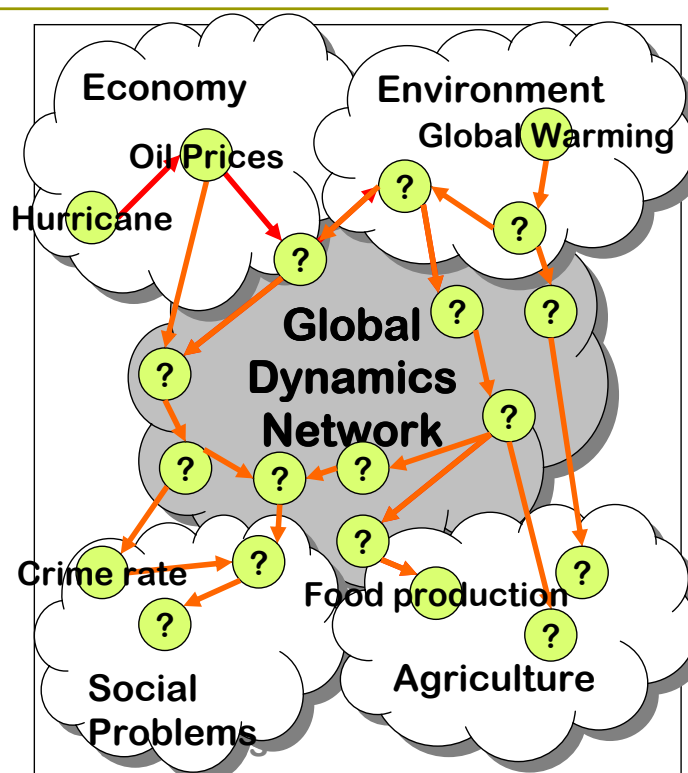
**Hideki Kawai[†], Kazuo Kunieda[†], Keiji Yamada[†],**
**Haruka Saito, Masaaki Tsuchida, Hironori Mizuguchi**
**[†] NEC C&C Innovation Research Laboratories**
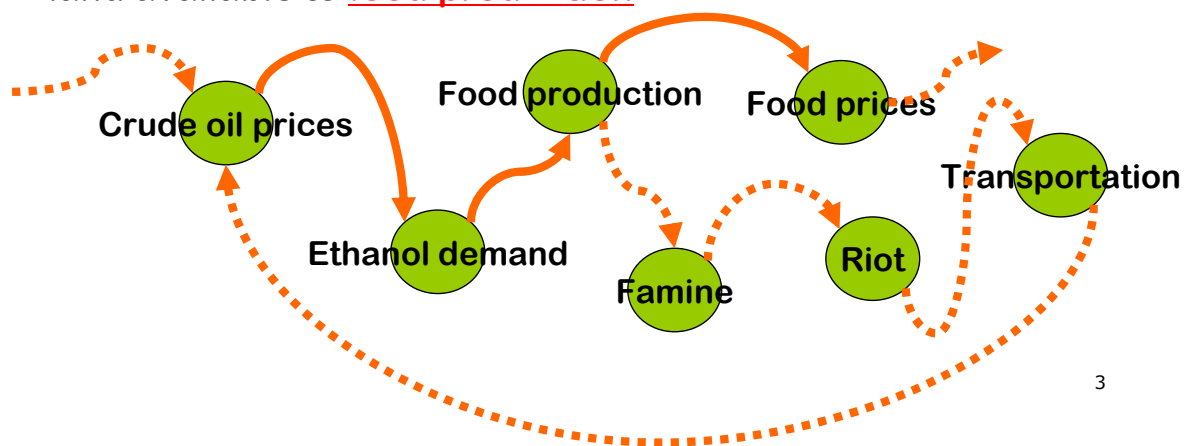**NEC Common Platform Software Research Laboratories**

1

# Background

- **Complex relations between various problems**
- **Causal relations**
  - **Butterfly effect**
  - **Ripple effect**
- **Global Dynamics**
  - **Global Solution**
  - **Idea Support**
- **Focusing on Statistical Terms**
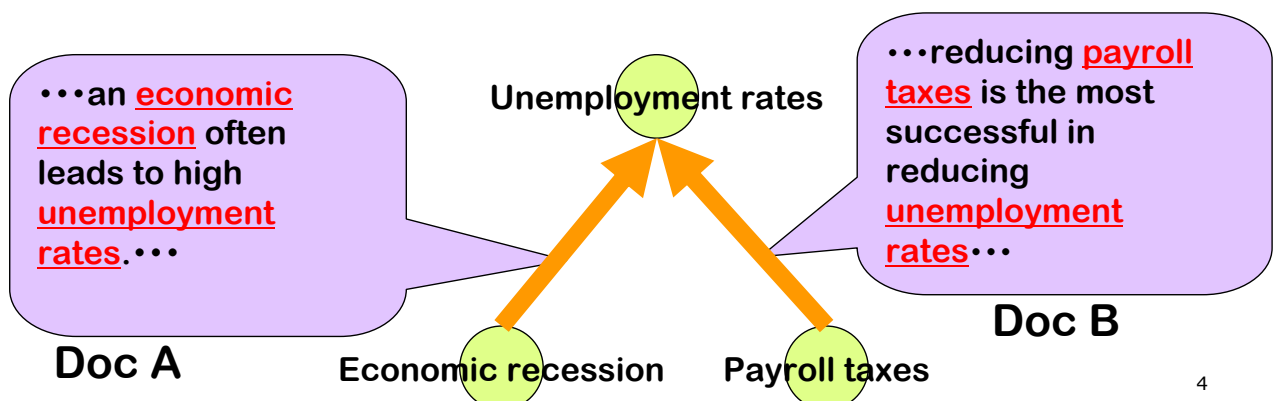
# Motivation

- **Web contains several billions of text pages**
- **Text contains information like causal relations**
  - Ex1: "the rise in <u>crude oil prices</u> stimulated <u>ethanol demand</u>"
  - Ex2: "farmers to answer to the <u>ethanol demand</u>, leave less land available to <u>food production</u>"

Crude oil prices — Ethanol demand — Food production — Food prices — Famine — Riot — Transportation

3

# Goal

- **Visualize Global Dynamics as a Statistical Term Network**
  - **Node：Statistical terms**
  - **Edge：Relationship between terms**

•••an <u>economic recession</u> often leads to high <u>unemployment rates</u>•••

Doc A

•••reducing <u>payroll taxes</u> is the most successful in reducing <u>unemployment rates</u>•••

Doc B

Unemployment rates

Economic recession    Payroll taxes
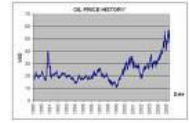
4

# What are we looking for?

- ☐ **Nodes:** <u>Expressions</u>
  - ■ <u>Statistical terms</u> (relative to a movement)
    - ☐ Ex. Unemployment rates, birth rates, crime rates, oil prices, corn prices
    - ☐ **Suffixes**: Rates, prices, costs
  - ■ <u>Events & Facts</u>
    - ☐ Ex. Hurricanes, riot, war
    - ☐ Ex. Urbanization, sustainable resources, global warning
    - ☐ **Classes**: natural disasters, social, economics

> ➢ **Extraction with the help of suffix & class patterns**

---

# What are we looking for?

- ☐ **Edges:** <u>Relations</u>
  - ■ **Link between a cause and its effect**
  - ■ **In text: Cause & effect are expressions linked via some verbs.**
  - ■ **For causal relation, in 80% of case: cause & effect are in the same sentence [1]**
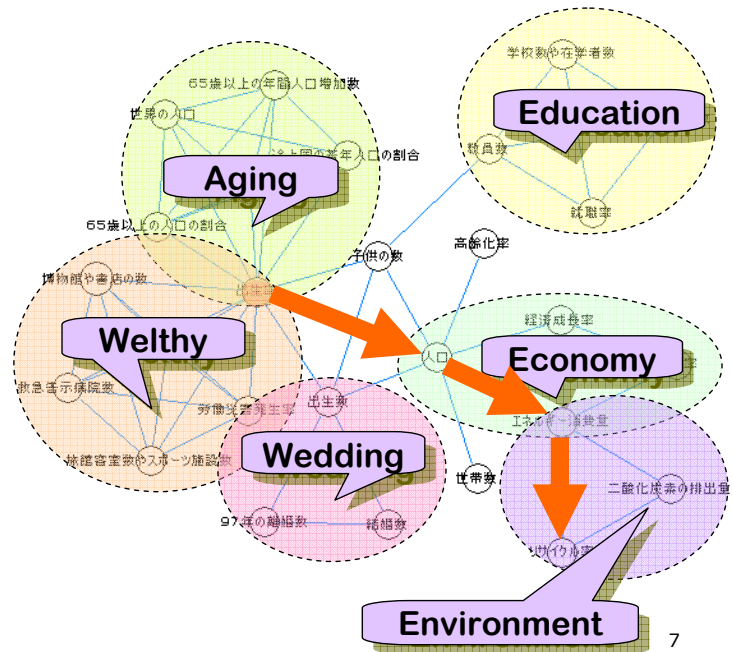
> ➢ **We will extract sentences containing a couple of expressions**

[1] Takashi Inui, *Creating an Annotated Corpus for the Analysis of Causal Relations*

# Previous Work [NTCIR-6, 2007]

- **Suffix-based Statistical term extraction**
  - 「**失業率** (unemployment rates)」
  - 「**原油価格** (oil prices)」
  - 「**物価指数** (consumer price index)」
- **Co-occurrence network of statistical terms**
  - Manually constructed

65歳以上の年齢人口増加数

世界の人口

Education

Aging

学校数や在学者数

教員数

65歳以上の人口の割合

試職率

博物館や書店の数

子供の数

高齢化率

経済成長率

Welthy

救急告示病院数

出生数

人口

Economy

労働災害発生率

エネルギー消費量

旅館客室数やスポーツ施設数

Wedding

世帯数

二酸化炭素の排出量

97年の離婚数

結婚数

リサイクル率

Environment

7

---

# Suffix-based Statistical Term Extraction

- 1) Find a suffix of statistical terms
  - ・・・ 増え て おり 、 ビール の 出荷 **台数** が ・・・ ○
  - ・・・ ある 場合 では テレビ の 生産 **台数** 競争 ・・・ ×

- 2) Scan leftward from the starting point to find the morpheme which is neither noun nor specific particles
  - ・・・ 増え て おり 、 ビール の 出荷 **台数** が ・・・
  - End ← Start

- 3) Extract morphemes between starting point to ending point as statistical terms
  - ・・・ 増え て おり 、 ビール の 出荷 **台数** が ・・・

**Statistical Term**

8

# Semantic Structure of Statistical Terms

## □ Base Form

- Shortest sequence of morphemes having a statistically valid meaning
  - □ **Unemployment rates, Oil prices**

---

**Various combinations:**
「unemployment rates」
「domestic unemployment rates」
「American unemployment rates in Mar. 1998」

## □ Modifiers

- **Object**
  - □ What is measured
    - **Beer、PC**
- **Subject**
  - □ Who measured
    - **Kirin、NEC**
- **Time Span**
  - □ When was it measured
    - **1999、Feb.**
- **Region**
  - □ Where was it measured
    - **Japan、America** 9
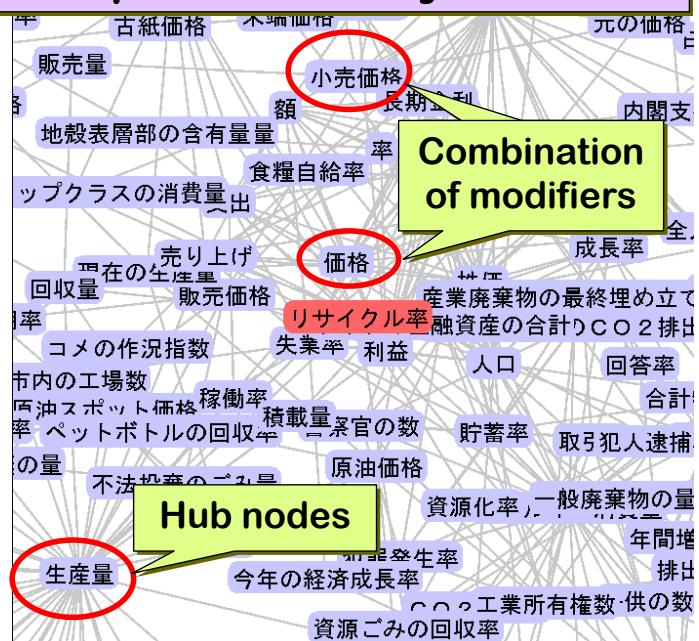
---

# Example of Statistical Term Network

## □ Too High Density

- **Difficult to see the relationship between statistical terms**

## □ Reasons

- **Hub nodes**
  - □ price、volume of production
- **Combination of modifiers**
  - □ price、retail price



2 hops from "recycle ratio"

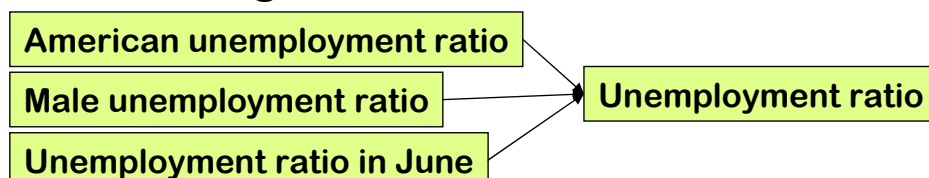Combination of modifiers

Hub nodes

10

# Challenges

- **Complexity of Network Structure**
  - **High dense network**
    - Clustering generates one big cluster
    - Threshold of co-occurrence does not work because most co-occurrence of statistical terms are only 1 or 2.

- **Complexity of Semantic Structure**
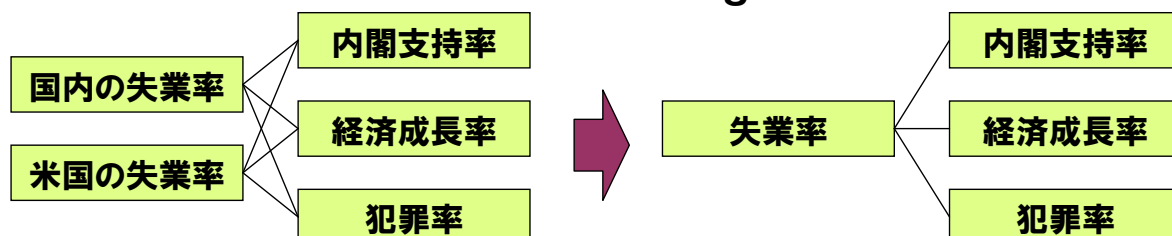  - **Appropriate combination of modifiers is not clear**
    - **Can we integrate?**

| American unemployment ratio |
| Male unemployment ratio | → | Unemployment ratio |
| Unemployment ratio in June |

11

# Our Approach

- **Simplify the Network Structure**
  - **Limit the degree of nodes**
    - Link only top $\omega$ terms
    - Main structure can be observed

- **Simplify the Semantic Structure**
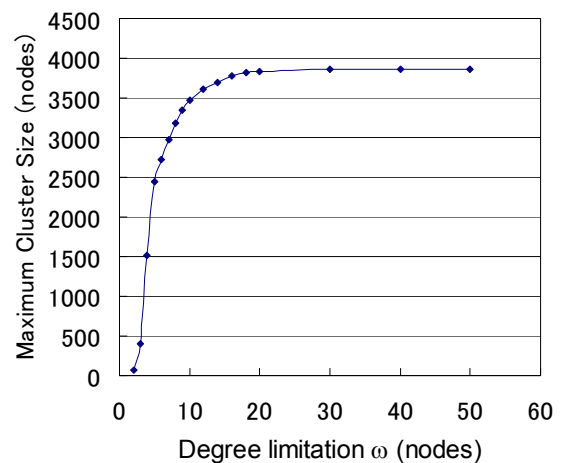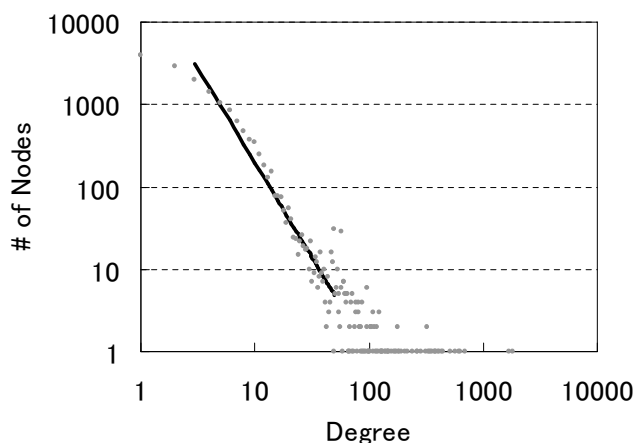  - **Integrate statistical terms which have common base form and co-occurring terms**

国内の失業率　米国の失業率 → 内閣支持率　経済成長率　犯罪率

⇒

失業率 → 内閣支持率　経済成長率　犯罪率

12

# Experimental Settings

- **Suffix dictionary**
  - **86 statistical terms tagged in MuST corpus**
- **Mainichi News 1998-1999 (Japanese)**
- **Extracted terms：8,600 words**
- **Degree parameter $\omega$**
  - **Investigated Maximum cluster size with $\omega$**
- **Visualization Tool**
  - **prefuse (http://prefuse.org/)**

13

# Degree Limitation Parameter $\omega$

- **Degree distribution of statistical terms follows power law**
- **Smaller Limitation Parameter $\omega$ divides the original network into small pieces**
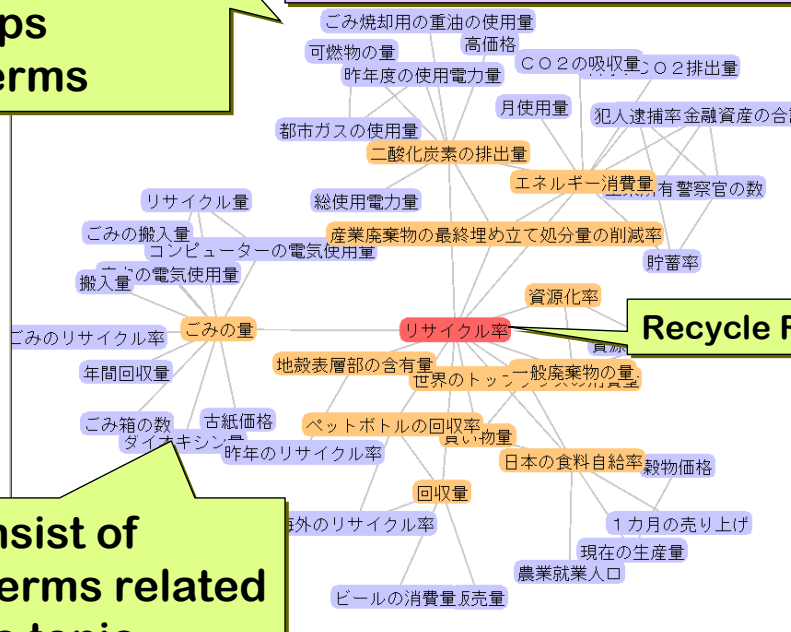- **We chose $\omega$ = 10 for the visualization.**



14

# Result of Semantic Structure Simplification

| Original Expression | Simplified Expression |
|---|---|
| わが国の温室効果ガスの総排出量<br>Domestic total greenhouse gas emission | 温室効果ガスの総排出量<br>Total greenhouse gas emission |
| 96年の温室効果ガスの総排出量<br>Total greenhouse gas emission in 1996 | |
| 12月のパソコン販売台数<br>PC sales volume in December | パソコン販売台数<br>PC sales volume |
| 秋葉原の電気街のパソコン販売台数<br>PC sales volume in Akihabara | |
| 埼玉県所沢市の野菜の価格<br>Vegetable prices in Tokorozawa city Saitama | 野菜の価格<br>Vegetable Prices |
| すべての野菜の価格<br>All vegetable prices | |

# Result of Network Structure Simplification

Easier to observe relationships between terms

2 hops from "Recycle Ratio"



Recycle Ratio

Cliques consist of statistical terms related to a specific topic

16

**7 hops from "Recycle Ratio"**

Recycle Ratio

17



Recycle Ratio
↓
Energy Consumption
↓
Amount of waste
↓
Dioxin
↓
Incidence of testicular cancer

Recycle Ratio
↓
Resource Recovery
↓
Sales of Beer
↓
Sales of PC
↓
Market Share

Recycle Ratio
↓
Unemployment Rates
↓
Economic Growth
↓
Interest rates
↓
Condominium Prices

Recycle Ratio

18

# Example: Energy Consumption



# Example: Amount of Waste

# Example: Amount of Dioxin

古紙価格
リサイクル率
ごみ箱の数
庁内の電気使用量
リサイクル量
ごみのリサイクル率

**Amount of waste**

ごみの量
コンピューターの電気使用量
年間回収量
し尿の量
人の数

ごみの搬入量

場合のダイオキシン摂取量
ダイオキシン排出量
食品別の摂取
許容量
分泌かく乱作用や体内への吸収率
日本人のダイオキシン摂取量

**Amount of dioxin**

ダイオキシン量
ダイオキシンの摂取量
アレルギーの発生率

**Incidence of allergy**
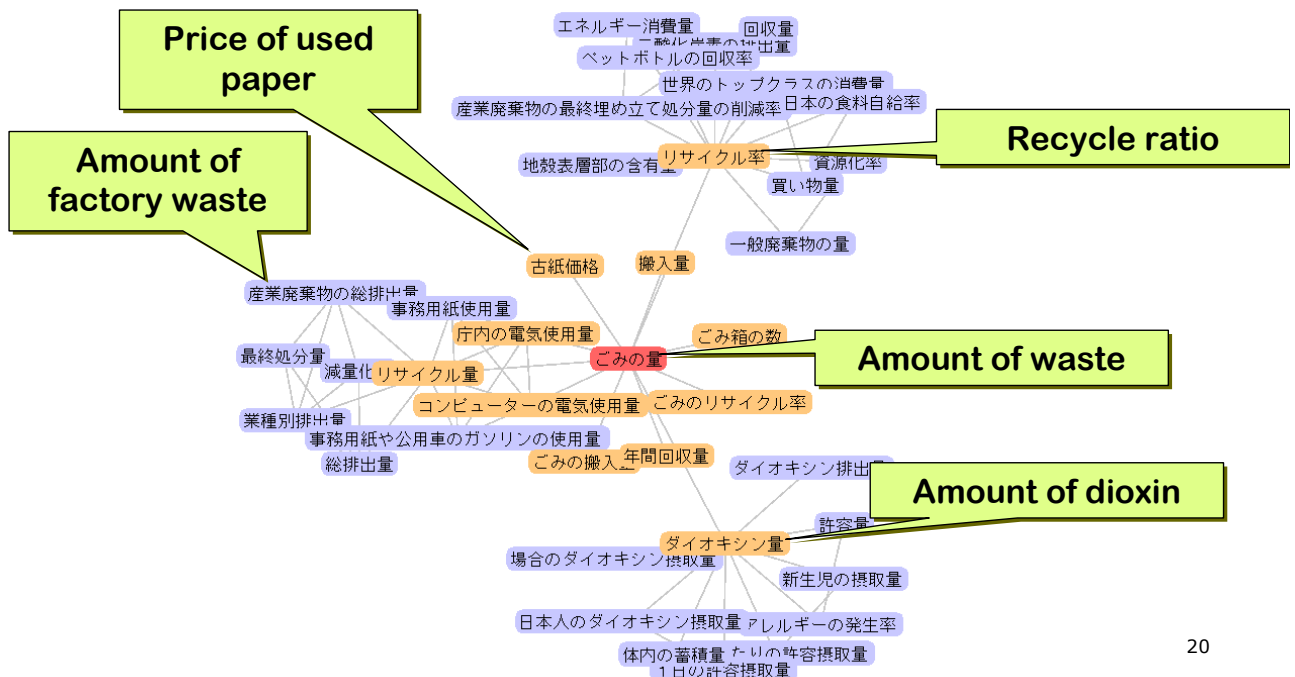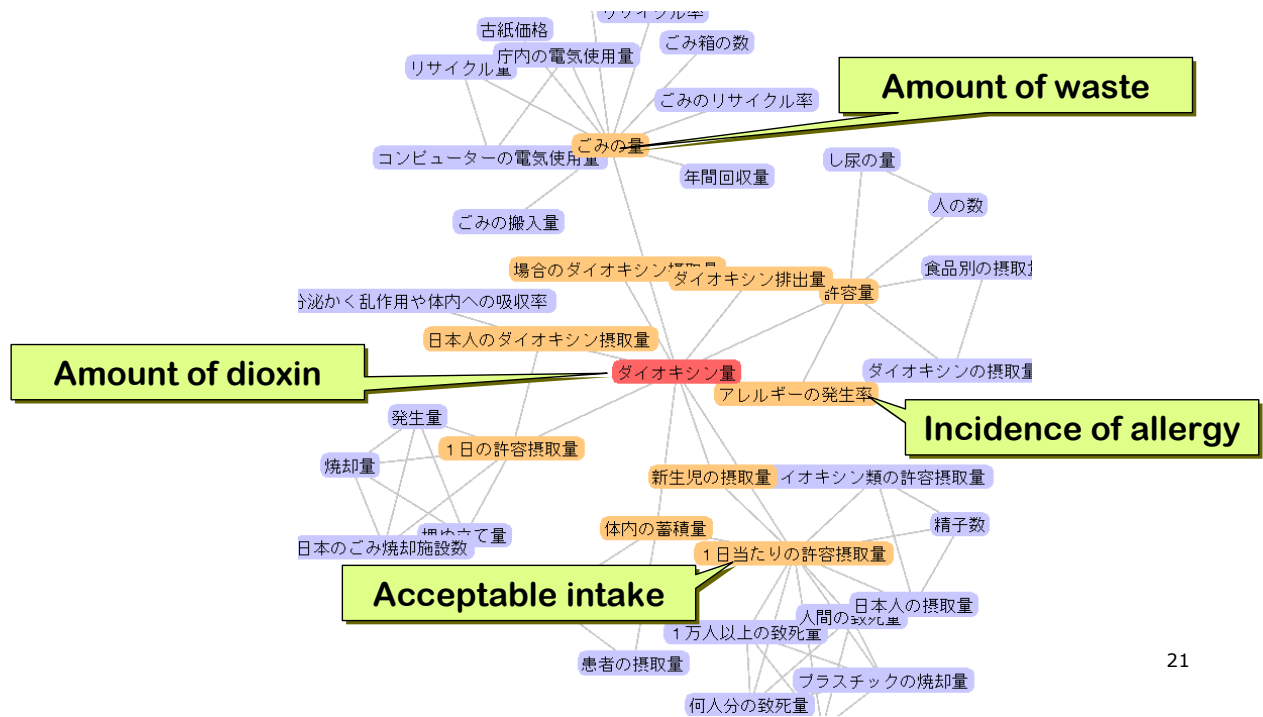
発生量
1日の許容摂取量
焼却量

新生児の摂取量　イオキシン類の許容摂取量
精子数

日本のごみ焼却施設数
埋め立て量
体内の蓄積量
1日当たりの許容摂取量

**Acceptable intake**

日本人の摂取量
人間の致死量
1万人以上の致死量
患者の摂取量
プラスチックの焼却量
何人分の致死量

21

# Examples: Unemployment Rates

お年玉の平均額
学生の就職内定率
高失業率
全体の求人倍率　ゴミ収集量
関連商品の売り上げ
倍率
サ　帯の平均消費性向

**Job openings**

雇用や収入
失業率や有効求人倍率
就職希望数
求人数

求人倍率
有効求職数
生活不安度指数
消費者態度指数　の全産業の経常利益
日米両国の完全失業率
GDP実質成

**Unemployment Rates**

全世帯消費支出　GDP成長率
完全失業率　完全失　完全失業者数
名目成長率

**Consumer spending**

昨年の完全失業者数
雇用者
完全失業率
米の完全失業率
就業者数
有効求人倍率

**Labor force population**

完全失業率
問題の完全失業率
手取り収入
同月の完全失業率
消費性向
対前年度比増減率
労働力人口
人の数
流産率
1世　消費支出　平均消費支
サラリーマン世帯
日本の完全失業率
サラリーマン世帯の平均消費支
過去最悪の失業率
許容量
消費者物価指数
日本人の配偶者のパート収入
とも初年度の支給合計額
均衡失業率　計額
免除　内の仮設住宅入居世帯数
し尿の量

22

# Examples: GDP Growth

**Land Prices and Stock Prices**

**House Prices**

**Unemployment Rates**

**GDP Growth**

**Consumer Spending**

**Share of the votes for Liberal Democratic Party**

株価や地価
銀行の不良債権 製造業の製品価格判断指数
株価純資産倍率
株価の下落率
当初計画値
ものの価格
住宅地価
地価や株価
中古マンション相場 住宅
6大都市市街地価格指数
銀行株価指数 東証株価指数
３大都市圏のマンション価格
住宅価格
新築マンショ
中古のマンション価格
購入価格
国内景気や収入
完全失業率
生活不安度指数 GDP成長率 住宅金融公
全世帯消費支出
同期の全産業の経常利益
自民党得票率 民党の支持率
成長率
名目成長率
政権への支持
票率
新設住宅着工戸数
民間消費支出 対前年増減率
政党費収入
政府支出
前年度比増減率
対前年度比増減率 パーティー収入 選挙での得票
卸売物価指数 政党の収入
騰落率
政治資金収入 収入額

23

# Example: Meal Size

**BMI**

**Cardiac rate**

**Exercise volume**

**Bone mass**

**Meal Size**

**Body fat percentage**

**Caloric intake**

エネルギー量 栄養所要量
中間の量
食事全体の量
体格指数
体脂肪の備蓄量 ＤＮＡ損傷の量
運動量
脂肪の量
心拍数
適度の量
食事量
体脂肪率
あごの幅の発 マッチ出場数
骨量
カロリー摂 家での食事の量
脂肪や内臓脂肪の量

24

# Conclusion

- **Visualization for Statistical Term Network**
  - **Simplification of Network Structure**
    - **Degree Limit Parameter $\omega$**
  - **Simplification of Semantic Structure**
    - **Integrate common base form and co-occurring terms**
  - **Experiment: Visualize a network consists of 8,600 statistical network**
- **Future Work**
  - **Exploit syntactic patterns about causal relation expression and make a direction on the statistical term network**
  - **Exploit Web corpus**

25