# Accelerating Biomolecular Nuclear Magnetic Resonance Assignment with A*

Joel Venzke, Paxten Johnson, Rachel Davis, John Emmons,
Katherine Roth, David Mascharka, Leah Robison,
Timothy Urness and Adina Kilpatrick

Department of Mathematics and Computer Science
Drake University

*joel.venzke@drake.edu*

April 10,2014

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick          Drake University

NMR Assignment with A*

**Introduction**
○
○○

NMR Assignment Background
○○

Automation Algorithm
○○○
○○○○○○○
○○○○

Conclusion
○○
○○○○

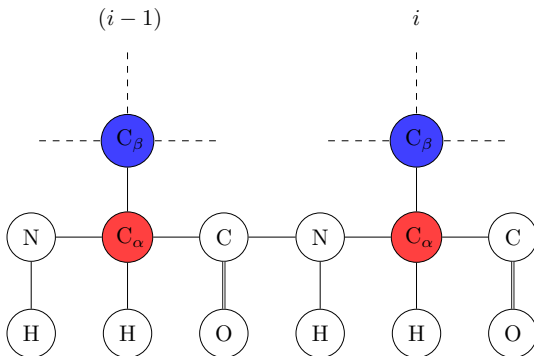# Overview

Motivation

# Motivation

- Nuclear Magnetic Resonance Spectroscopy
  - Gain knowledge about protein structure
  - Study how mutations lead to diseases
- Problems
  - Generates large amounts of data
  - Data analysis is slow and error prone
- Goal
  - Automate the assignment process
  - Decrease human error
  - Increase productivity

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick    Drake University

NMR Assignment with A*

# Nuclear Magnetic Resonance (NMR)

- Used to obtain structural information
  - Chemical shift values
- HNCACB experiment
  - Generates $C_\alpha$ and $C_\beta$ residue $i$ and $i-1$
- CBCA(CO) NH experiment
  - Generates $C_\alpha$ and $C_\beta$ for residue $i$
  - Confirms residue data

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick        Drake University

NMR Assignment with A*

Introduction
○
○●

NMR Assignment Background
○○

Automation Algorithm
○○○
○○○○○○○
○○○○

Conclusion
○○
○○○○

Nuclear Magnetic Resonance Spectroscopy

# Chemical Shift Values



HNCACB

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick          Drake University

NMR Assignment with A*

| Introduction | NMR Assignment Background | Automation Algorithm | Conclusion |
|---|---|---|---|
| ○ | ●○ | ○○○ | ○○ |
| ○○ | | ○○○○○○○ | ○○○○ |
| | | ○○○○ | |

Data Collection and Manual Assignment

# Manual Methods

- Most time consuming part
- Missing and ambiguous data forces chunks to be skipped
- Prone to human error

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick    Drake University

NMR Assignment with A*

Introduction
○
○○

NMR Assignment Background
○●

Automation Algorithm
○○○
○○○○○○○
○○○○

Conclusion
○○
○○○○

Data Collection and Manual Assignment

# Timeline

Protein
Production
at least 5 days

Data Assignment
20 days to 9 months

NMR
Experiments
1-2 days

[1]

Preprocessing

# Automating Assingment

- Initialization
- Generating child nodes
- Goal State
- Solution State

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick      Drake University

NMR Assignment with A*

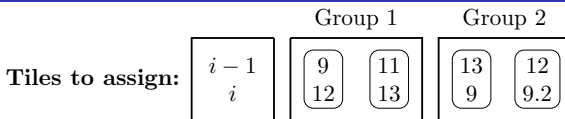| Introduction | NMR Assignment Background | Automation Algorithm | Conclusion |
|---|---|---|---|
| ○ | ○○ | ○●○ | ○○ |
| ○○ | | ○○○○○○○ | ○○○○ |
| | | ○○○○ | |

Preprocessing

# Initialization

- Expected amino acid sequence
  - Converted to expected chemical shift values
  - Stored as the reference protein chain
- NMR experiment's chemical shift data
  - $C_\alpha$ and $C_\beta$ for residue $i$ and $i-1$
  - Stored in a tile
- Missing data
  - Place holder tile generation
- Grouping

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick     Drake University

NMR Assignment with A*

# Grouping



[2]

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick    Drake University

NMR Assignment with A*

Assignment

# Starting the assignment



|                |       | Group 1 | Group 2 |
|----------------|-------|---------|---------|
| **Tiles to assign:** | $\begin{matrix} i-1 \\ i \end{matrix}$ | $\begin{matrix} 9 \\ 12 \end{matrix}$ $\begin{matrix} 11 \\ 13 \end{matrix}$ | $\begin{matrix} 13 \\ 9 \end{matrix}$ $\begin{matrix} 12 \\ 9.2 \end{matrix}$ |

**Reference Protein Chain**                 **Nodes**

| Chemical Shift | Group |
|----------------|-------|
| 13.5           | 1     |
| 9.5            | 2     |
| 11.4           | 1     |
| 8.8            | 2     |

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick    Drake University

NMR Assignment with A*

Introduction
○
○○

NMR Assignment Background
○○

Automation Algorithm
○○○
○●○○○○○
○○○○

Conclusion
○○
○○○○

Assignment

# Starting the assignment

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick    Drake University

NMR Assignment with A*

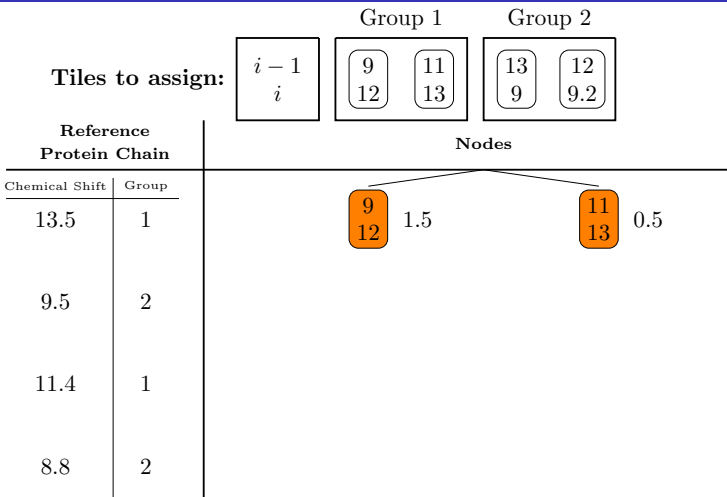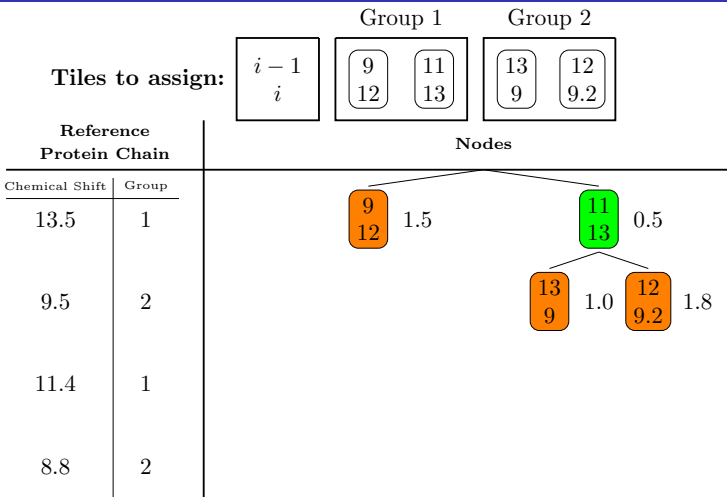| Introduction | NMR Assignment Background | Automation Algorithm | Conclusion |
|---|---|---|---|
| ○ | ○○ | ○○○ | ○○ |
| ○○ | | ○○●○○○○ | ○○○○ |
| | | ○○○○ | |

Assignment

# Cost Calculation

- Accuracy matching the protein chain residue
- Accuracy matching the tile above current tile
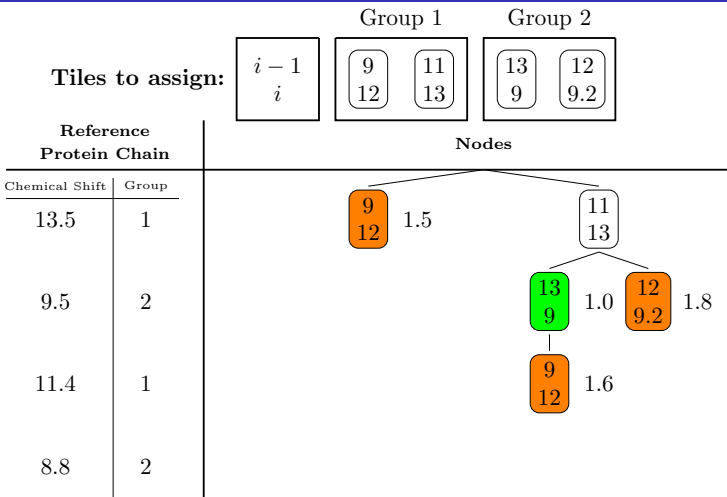- Cost of placing all previous tiles

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick     Drake University

NMR Assignment with A*

Introduction
○
○○

NMR Assignment Background
○○

Automation Algorithm
○○○
○○○●○○○
○○○○

Conclusion
○○
○○○○

Assignment

# Generating child nodes

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick    Drake University

NMR Assignment with A*

Introduction
○
○○

NMR Assignment Background
○○

Automation Algorithm
○○○
○○○○●○○
○○○○

Conclusion
○○
○○○○

Assignment

# Generating child nodes

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick          Drake University

NMR Assignment with A*

Introduction
○
○○

NMR Assignment Background
○○

Automation Algorithm
○○○
○○○○○●○
○○○○

Conclusion
○○
○○○○

Assignment

# Generating child nodes

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick    Drake University

NMR Assignment with A*

Introduction
○
○○

NMR Assignment Background
○○

Automation Algorithm
○○○
○○○○○○●
○○○○

Conclusion
○○
○○○○

Assignment

# Generating child nodes

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick          Drake University

NMR Assignment with A*

| Introduction | NMR Assignment Background | Automation Algorithm | Conclusion |
|---|---|---|---|
| ○ | ○○ | ○○○ | ○○ |
| ○○ | | ○○○○○○○ | ○○○○ |
| | | ●○○○ | |

Goal State

# Goal State

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick          Drake University

NMR Assignment with A*

Introduction
○
○○

NMR Assignment Background
○○

Automation Algorithm
○○○
○○○○○○○
○●○○

Conclusion
○○
○○○○

Goal State

# Goal State

Introduction
○
○○

NMR Assignment Background
○○

Automation Algorithm
○○○
○○○○○○○
○○●○

Conclusion
○○
○○○○

Goal State

# Goal State

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick    Drake University

NMR Assignment with A*

| Introduction | NMR Assignment Background | Automation Algorithm | Conclusion |
| ○ | ○○ | ○○○ | ○○ |
| ○○ | | ○○○○○○○ | ○○○○ |
| | | ○○○● | |

Goal State

# Solution State

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick    Drake University

NMR Assignment with A*

Introduction
○
○○

NMR Assignment Background
○○

Automation Algorithm
○○○
○○○○○○○
○○○○

Conclusion
●○
○○○○

Results

# Time of Assignment

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick          Drake University

NMR Assignment with A*

Introduction
○
○○

NMR Assignment Background
○○

Automation Algorithm
○○○
○○○○○○○
○○○○

Conclusion
○●
○○○○

Results

# Child Nodes Generated

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick     Drake University

NMR Assignment with A*

# Future Goals

- Parallelization
  - Decrease assignment time
  - Allow for larger data sets
- Machine learning
  - Optimize cost calculation
  - Increase accuracy of assignment

# Acknowledgments

- Dr. Tim Urness (Mathematics and Computer Science)
- Dr. Adina Kilpatrick (Physics)
- Rachel Davis (research colleague)
- John Emmons (research colleague)
- Katherine Roth (research colleague)
- David Mascharka (research colleague)
- Leah Robison (research colleague)

| Introduction | NMR Assignment Background | Automation Algorithm | Conclusion |
|---|---|---|---|
| ○ | ○○ | ○○○ | ○○ |
| ○○ | | ○○○○○○○ | ○○●○ |
| | | ○○○○ | |

Outlook

# Bibliography

📄 Babak Alipanahi, Xin Gao, Emre Karakoc, Frank Balbach, Shuai Cheng Li, Guangyu Feng, Logan Donaldson and Ming Li, *Error tolerant NMR backbone resonance assignment and automated structure generation.*, Journal of bioinformatics and computational biology, **9** (2011), 15–41.

📄 Sean Cahill and Mark Girvin.
*Introduction to 3d triple resonance experiments.*
2012.

📄 Peter Guntert, *Automated structure determination from NMR spectra*, European Biophysics Journal, **38** (2009), 129–143.

📄 Flemming M. Poulsen.
*A brief introduction to nmr spectroscopy of proteins.*

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick     Drake University

NMR Assignment with A*

| Introduction | NMR Assignment Background | Automation Algorithm | Conclusion |
|---|---|---|---|
| ○ | ○○ | ○○○ | ○○ |
| ○○ | | ○○○○○○○ | ○○○● |
| | | ○○○○ | |

Outlook

# Thank You

J. Venzke, P. Johnson, R. Davis, J. Emmons, K. Roth, D. Mascharka, L. Robison, T. Urness, A. Kilpatrick      Drake University

NMR Assignment with A*