# Improving Centruflow Using Semantic Web Technologies

A thesis presented in partial fulfillment of the requirements for
the degree of Master of Science in Computer Science at
Massey University, Palmerston North, New Zealand.

Jonathan Andrew Giles

2007

# Abstract

Centruflow is an application that can be used to visualise structured data. It does this by drawing graphs, allowing for users to explore information relationships that may not be visible or easily understood otherwise. This helps users to gain a better understanding of their organisation and to communicate more effectively. In earlier versions of Centruflow, it was difficult to develop new functionality as it was built using a relatively unsupported and proprietary visualisation toolkit. In addition, there were major issues surrounding information currency and trust. Something had to be done, and this was a sub-project of this thesis.

The main purpose of this thesis however was to research and develop a set of mathematical algorithms to infer implicit relationships in Centruflow data sources. Once these implicit relationships were found, we could make them explicit by showing them within Centruflow. To enable this, relationships were to be calculated based on providing users with the ability to 'tag' resources with metadata. We believed that by using this tagging metadata, Centruflow could offer users far more insight into their own data.

Implementing this was not a straight-forward task, as it required a considerable amount of research and development to be undertaken to understand and appreciate technologies that could help us in our goal. Our focus was primarily on technologies and approaches common in the semantic web and 'Web 2.0' areas. By pursuing semantic web technologies, we ensured that Centruflow would be considerably more standards-compliant than it was previously. At the conclusion of our development period, Centruflow had been rather substantially 'retrofitted', with all proprietary technologies replaced with equivalent semantic web technologies. The result of this is that Centruflow is now positioned on the forefront of the semantic web wave, allowing for far more comprehensive and rapid visualisation of a far larger set of readily-available data than what was possible previously.

Having implemented all necessary functionality, we validated our approach and were pleased to find that our improvements led to a considerably more intelligent and useful Centruflow application than was previously available. This functionality is now available as part of 'Centruflow 3.0', which will be publicly released in March 2008. Finally, we conclude this thesis with a discussion on the future work that should be undertaken to improve on the current release.

# Acknowledgements

# Contents

x

# List of Figures

# List of Tables

# List of Listings