

Data Management Solutions

nf-core hackathon

July 15 2020

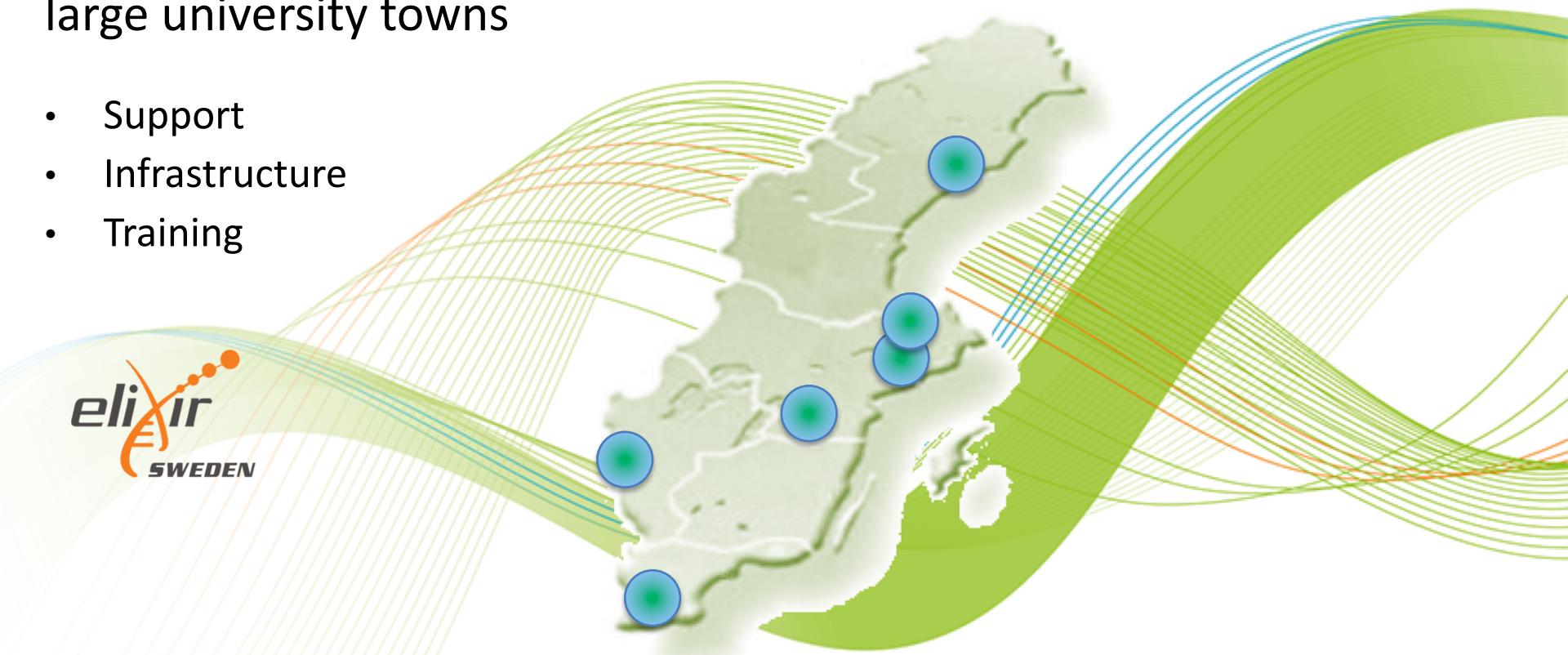
Elin Kronander

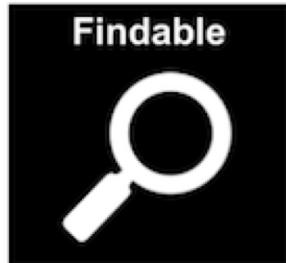
elin.kronander@nbis.se

Bioinformatics Platform

Distributed research infrastructure with nodes at each of the 6 large university towns

- Support
- Infrastructure
- Training





To be useful for others data should be **FAIR**
... for both Machines and Humans

www.nature.com/scientificdata

SCIENTIFIC DATA



OPEN

SUBJECT CATEGORIES

- » Research data
- » Publication characteristics

Received: 10 December 2015
Accepted: 12 February 2016
Published: 15 March 2016

Comment: The FAIR Guiding Principles for scientific data management and stewardship

Mark D. Wilkinson et al.*

There is an urgent need to improve the infrastructure supporting the reuse of scholarly data. A diverse set of stakeholders—representing academia, industry, funding agencies, and scholarly publishers—have come together to design and jointly endorse a concise and measurable set of principles that we refer to as the FAIR Data Principles. The intent is that these may act as a guideline for those wishing to enhance the reusability of their data holdings. Distinct from peer initiatives that focus on the human scholar, the FAIR Principles put specific emphasis on enhancing the ability of machines to automatically find and use the data, in addition to supporting its reuse by individuals. This Comment is the first formal publication of the FAIR Principles, and includes the rationale behind them, and some exemplary implementations in the community.

Supporting discovery through good data management
Good data management is not a goal in itself, but rather is the key conduit leading to knowledge discovery and innovation, and to subsequent data and knowledge integration and reuse by the community after the data publication process. Unfortunately, the existing digital ecosystem surrounding scholarly data publication prevents us from extracting maximum benefit from our research investments (e.g., ref. 1). Partially in response to this, science funders, publishers and

Wilkinson, Mark et al. "The FAIR Guiding Principles for scientific data management and stewardship".
Scientific Data 3, Article number: 160018 (2016)
<http://dx.doi.org/10.1038/sdata.2016.18>

FAIR Principles



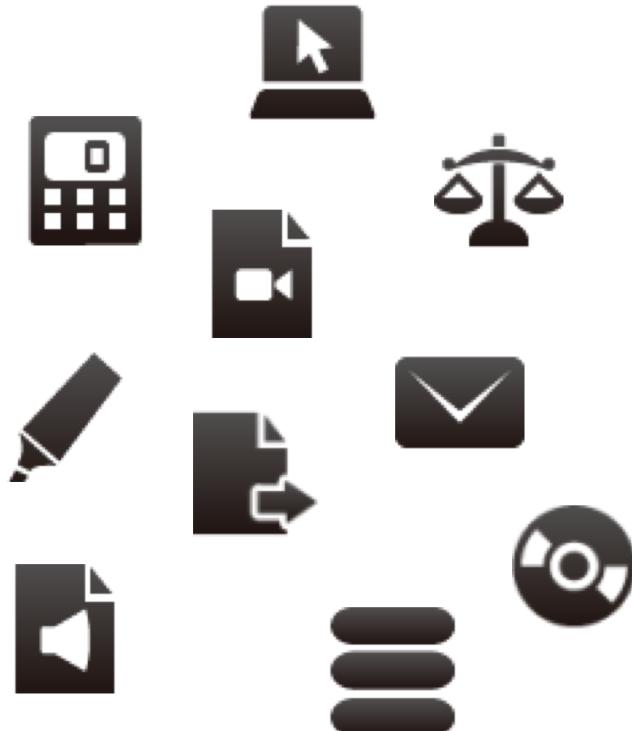


dbSNP
Short Genetic Variations

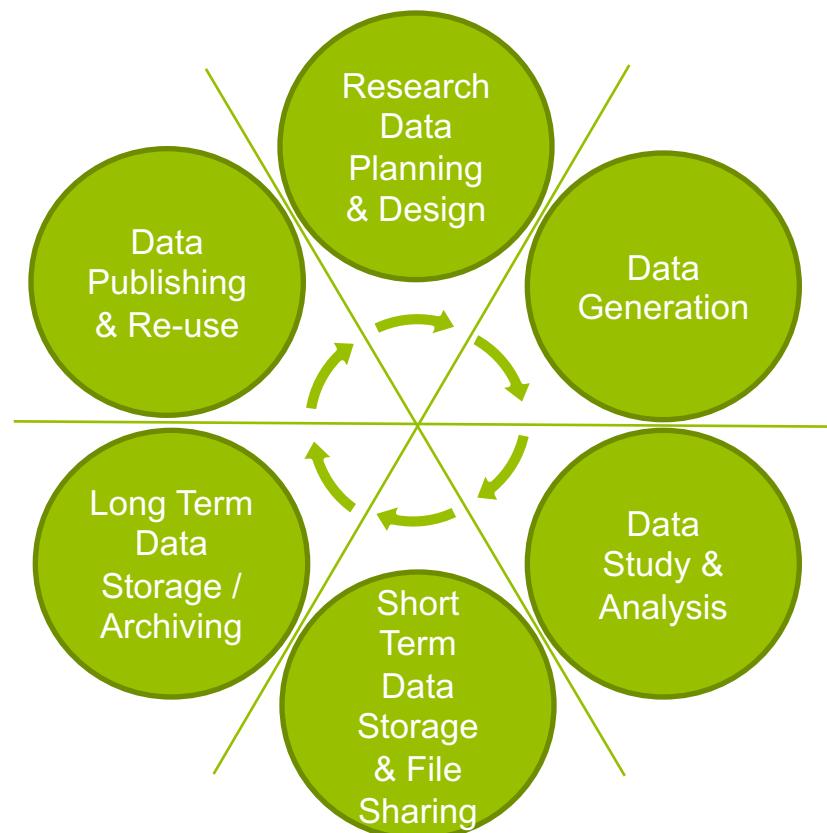


- Good way to make data **FAIR**
- Domain-specific metadata standards

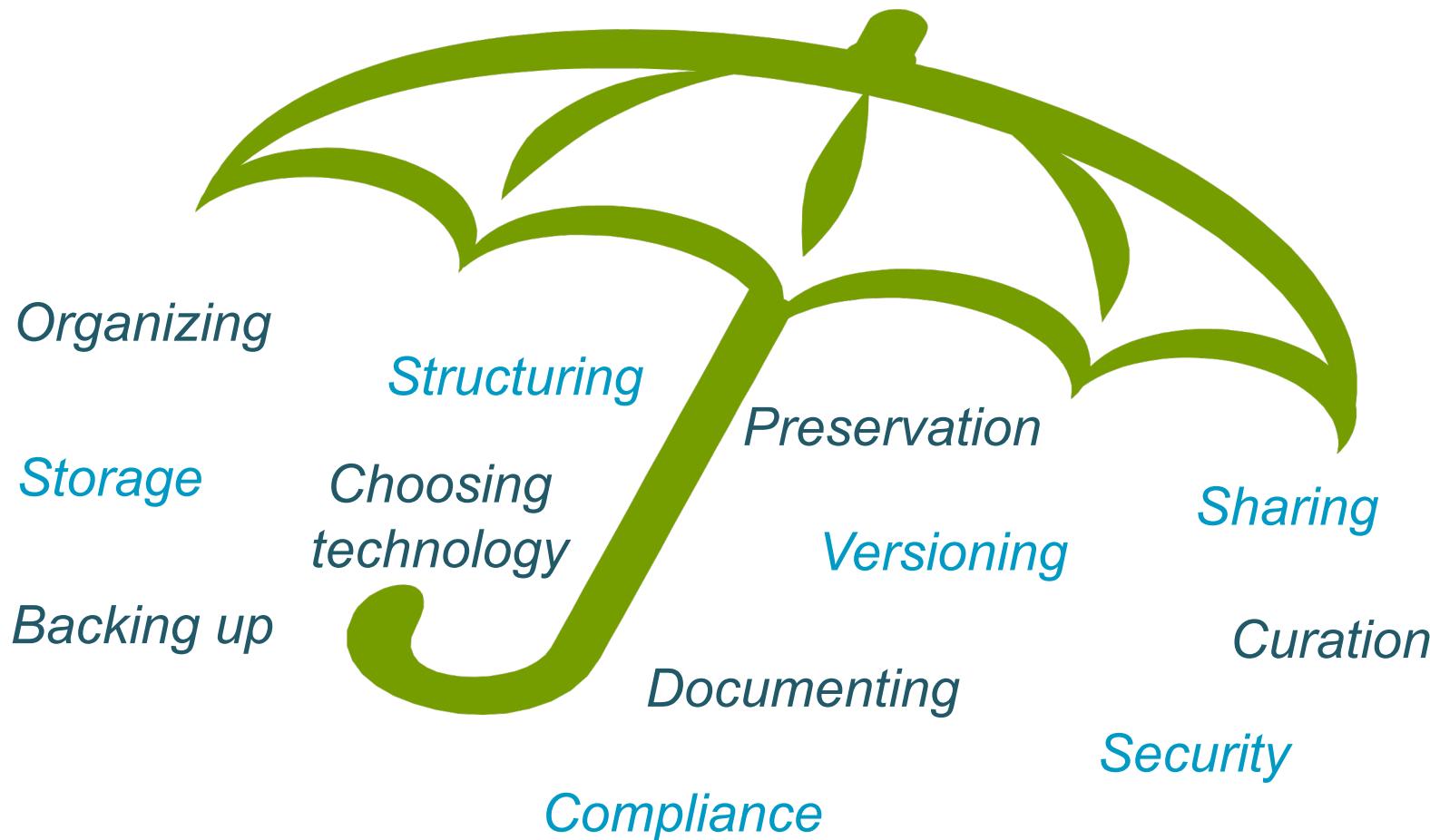
Any information you
use in your research



The Research Data Life Cycle



Purpose: Making the research process as efficient as possible



- Revisable Document
- New & existing data
- Throughout data life cycle

1. Description of data

- What type of data will be generated
- From which type of samples will the data be generated
- What additional data will be collected

2. Documentation and data quality

- How will the analysis be documented
- Which metadata standards will be used
- What quality measures will be used
- Which file formats will be used
- What is your strategy concerning versioning

3. Storage and Backup

- How are the data stored?
- Are there back-up systems
- What is the estimated size of data
- Do you need to restrict access to data

4. Legal and ethical aspects

- Is sensitive human data part of your project
- Are there agreements in place with other stakeholders

5. Accessibility and long-term storage

- How and where will the data be shared
- How and where will the data be stored after the project's completion

DMP tools

DMPonline

 Home Public DMPs Funder requirements Help Language ▾

Welcome

DMPonline helps you to create, review, and share data management plans that meet institutional and funder requirements. It is provided by the Digital Curation Centre (DCC).

Join the growing international community that have adopted DMPonline:

 17,622 Users	 203 Organisations
 23,083 Plans	 89 Countries

Some funders mandate the use of DMPonline, while others point to it as a useful option. You can [download funder templates](#) without logging in, but the tool provides tailored guidance and example answers from the DCC and many research organisations. Why not sign up for an account and try it out?

 My Dashboard Create plans Reference ▾ Help Language ▾ Rob Hooft ▾

DMP for a ZonMw Project

Project Details Plan overview Data Section Enabling Technologies Hotels Data management ZonMw Share Downloaded

expand all [collapse all] 0/29 answered

- 1. General information (0 / 11) +
- 2. Legislation and regulations (0 / 2) +
- 3. Findable (0 / 4) +
- 4. Accessible (0 / 3) +
- 5. Interoperable (0 / 4) +
- 6. Reusable (0 / 0) +
- 7. Sustainable data storage (0 / 5) +

<https://dmponline.dcc.ac.uk/>

ELIXIR Data Stewardship Wizard

Go to App


DSW
 DATA STEWARDSHIP WIZARD

Smart Data Management Plans for FAIR Open Science
 For Serious Researchers and Data Stewards

SciLifeLab DSW

Current Phase Before Submitting the Proposal

IV. Data storage and backup

1 What is the estimated total size of the data?

The (sequencing) facility should be able to tell you roughly how much space the raw sample(s) will take. When you're working with the data, it usually expands by a factor ranging between 50%-300%, and you will need to account for this.

Desirable: Before Submitting the DMP

a. Less than 1 TB
 b. Between 1 TB and 10 TB
 c. Between 10 TB and 50 TB
 d. Between 50 TB and 100 TB
 e. More than 100 TB

2 Where will the data be stored during the research

Consider submitting the read data to ENA early in the project as a back-up copy of the "raw" data. (Note that human data should not be submitted to the ENA)

Desirable: Before Submitting the DMP
 External Links: ENA

a. SNIC center

<https://ds-wizard.org/>

<https://dsw.scilifelab.se>

Home Pipelines Tools Usage Developers About [Join nf-core](#)

Data Management

How to plan your project, estimate resources, and share your results.

Getting started
Installation
Pipeline configuration
Running offline
nf-core tutorial
Reference genomes
Data Management
Troubleshooting
Nextflow resources

[Edit](#)

Data management

Funding agencies are recognizing the importance of research data management and some now request detailed Data Management Plans (DMP) as part of the grant application. Research data management concerns the organization, storage, preservation, and sharing of data that is collected or analyzed during a research project. Proper planning and data management facilitates sharing and allows others to validate and reuse the data. Guidance is provided below to aid the creation of DMPs, estimate resources needed by nf-core workflows, and how to share the resulting data.

Data Management Plan

A Data Management Plan (DMP) is a revisable document explaining how you intend to handle new and existing data, during and following the conclusion of your research project. It is wise to write a DMP as early as possible, using either a tool provided by your host institution or for example [DS Wizard](#) or [DMP Online](#). Ethical and legal considerations regarding the data will depend on where the research is conducted, this is especially true for projects including sensitive human data. For more information about the Swedish context, please review this page on [Sensitive personal data](#).

Data storage and computational resources

To estimate computational resources needed for a specific pipeline please see the pipeline documentation. This lists the different output filetypes you can expect. In the future we hope to automate a full run of each pipeline after every release. The pipeline docs will then show a full set of results from a real run, along with all file sizes. This can then be used as a guide as to what to expect for your

Data Sharing exercise

1. Select a dataset or a project you are working on



2. Identify a suitable repository



3. Check the repository guidelines for:

- file formats
 - raw data
 - processed data
- metadata
- recommended vocabularies
- quality measures
- licensing



4. Prepare your data to conform to the repository guidelines



5. Submit your raw data



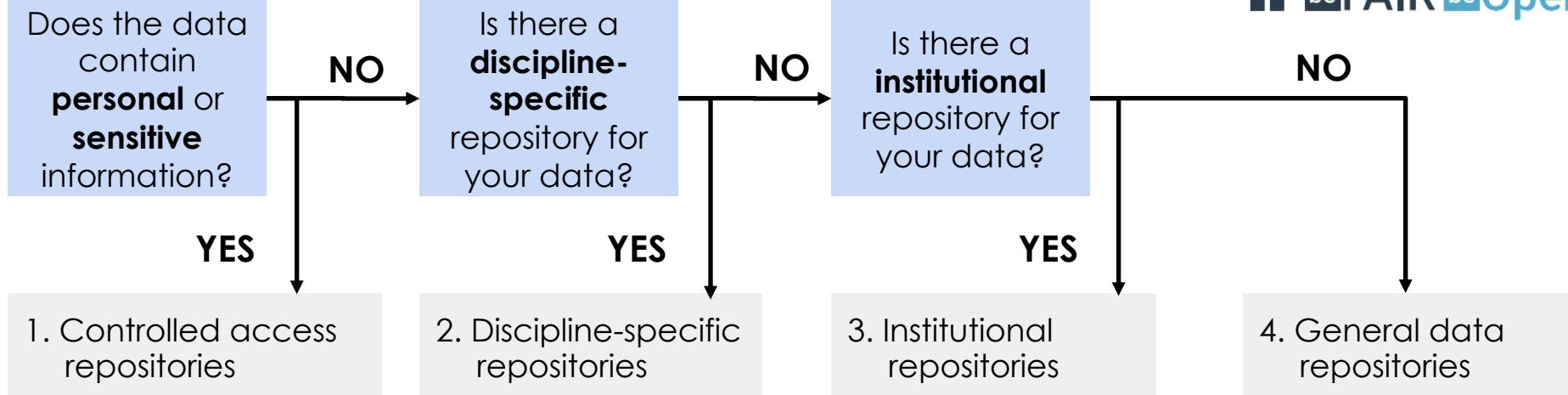
- easier to provide metadata
- ready for publication
- backup in different location
- increase citations

- processed data can be linked later
- use embargo if necessary

Choosing a repository

F1000

 beFAIR  beOpen



dbSNP
Short Genetic Variations



Etc...



The following metadata fields are supported in the manifest file:

STUDY: Study accession or unique name (alias)

SAMPLE: Sample accession or unique name (alias)

NAME: Unique experiment name

PLATFORM: [See permitted values](#). Not needed if INSTRUMENT is provided.

INSTRUMENT: [See permitted values](#)

INSERT_SIZE: Insert size for paired reads

LIBRARY_NAME: Library name (optional)

LIBRARY_SOURCE: [See permitted values](#)

LIBRARY_SELECTION: [See permitted values](#)

LIBRARY_STRATEGY: [See permitted values](#)

DESCRIPTION: free text library description (optional)

Permitted values for library source

GENOMIC: Genomic DNA (includes PCR products from genomic DNA).

GENOMIC SINGLE CELL:

TRANSCRIPTOMIC: Transcription products or non genomic DNA (EST, cDNA, RT-PCR, screened libraries).

TRANSCRIPTOMIC SINGLE CELL:

METAGENOMIC: Mixed material from metagenome.

METATRANSCRIPTOMIC: Transcription products from community targets

SYNTHETIC: Synthetic DNA.

VIRAL RNA: Viral RNA.

OTHER: Other, unspecified, or unknown library source material.

- Consider doing a Data Management Plan for your project
 - How do you ensure that your research output is FAIR?
- Plan for submitting “raw data” to public repositories as early as possible
- Organize project metadata from the start
 - In ways that makes it easy to submit to public repositories
 - Use available standards
- Pick a thought-through file and folder structure organization for your computational analyses
- Strive for reproducibility
 - Data & Code
- Be aware that there are legal aspects to processing human data
- *Ask for help if you need it!*

*“Your primary
collaborator is
yourself six months
from now, and your
past self doesn’t
answer e-mails,”*

-Rachael Ainsworth