Why you cannot "apt install" an LLM (yet)

How Debian is rethinking its guidelines for AI SFSCON 2025

About me

- Andrea Pappacoda <tachi>
- Debian contributor since 2021
- Debian Developer since 2024
- Computer science student specializing in cybersecurity (sometimes)

Intro

So, why cannot you apt-install an LLM?

It's a bit of a mess!

- In the value chain
- Things change with time
- We need free datasets

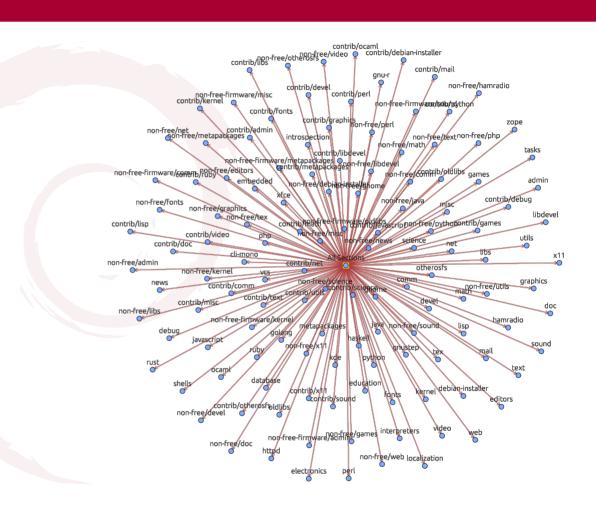
Debian is many things

• GNU/Linux "distro"

Debian is many things

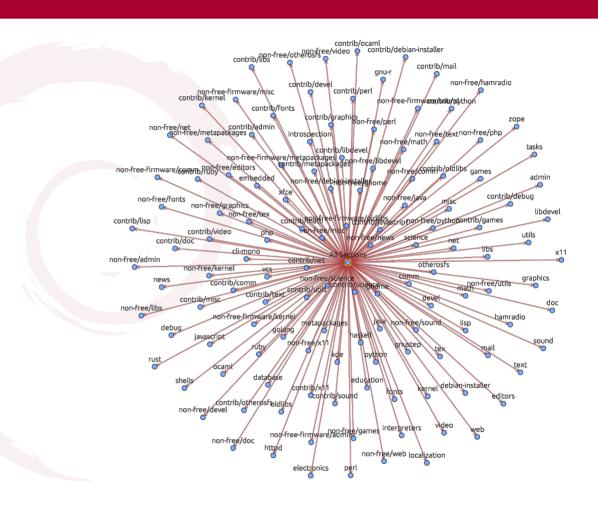
- GNU/Linux "distro"
- Distribution of Free Software

We distribute a lot of software.



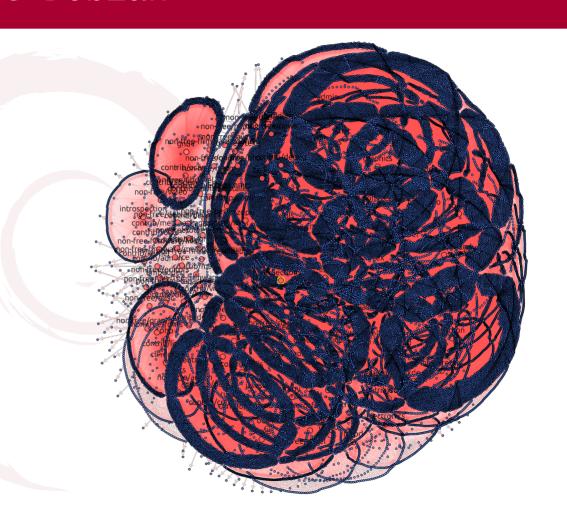
We distribute a lot of software.

If every dot is a category... how many packages are there?



...more than you could imagine

...more than you could imagine



It's a community of people



Birth of a package

- We find something interesting
- We study it a bit
- We tweak it if needed
- We build it
- We wrap it in a package

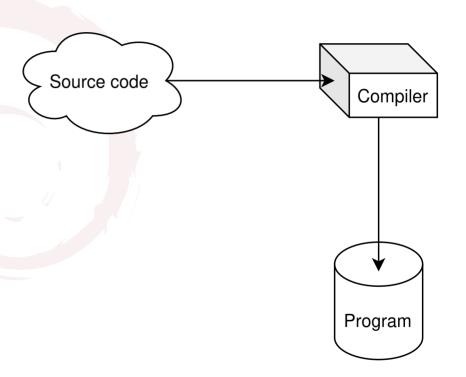
Birth of a package

- We find something interesting
- We study it a bit
- We tweak it if needed
- We build it
- We wrap it in a package

Birth of a package

Building a software entails

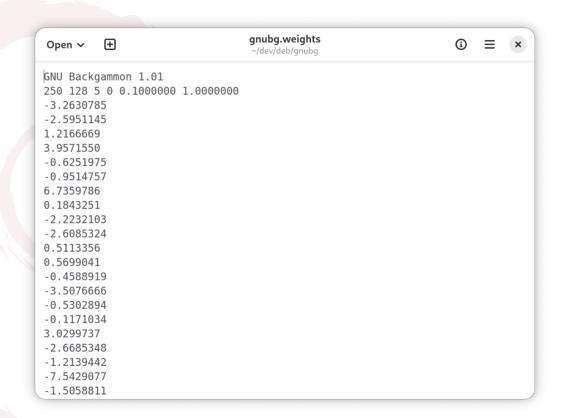
- The source code
- A compiler
- The resulting program



- A backgammon game and engine
- In Debian since 2001
- Can be used to play against the computer



- How can the computer play backgammon?
- What are all those numbers?
- Did someone write them all?



- For more than 20 years, nobody cared
- With LLMs becoming more relevant, old practice is put under discussion

- For more than 20 years, nobody cared
- With LLMs becoming more relevant, old practice is put under discussion

gnubg, at least, comes with neural network weights that do not have source code under this definition. I have to admit that while I was maintaining the package I didn't give this a ton of a thought because it predates the whole LLM craze.

- -

Russ Alberry <rra@debian.org>

- For more than 20 years, nobody cared
- With LLMs becoming more relevant, old practice is put under discussion

gnubg, at least, comes with neural network weights that do not have source code under this definition. I have to admit that while I was maintaining the package I didn't give this a ton of a thought because it predates the whole LLM craze.

_ _

Russ Alberry <rra@debian.org>

- Boring way of spelling out our goals
- What's "source code" in a neural network?

Debian Social Contract

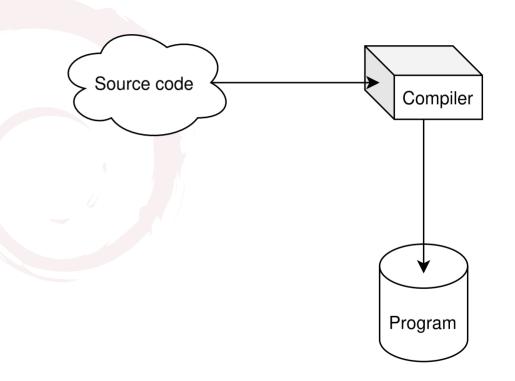
Version 1.2 ratified on October 1st, 2022.

Supersedes <u>Version 1.1</u> ratified on April 26th, 2004, and <u>Version 1.0</u> ratified on July 5, 1997.

Debian, the producers of the Debian system, have created the **Debian Social Contract**. The <u>Debian Free Software Guidelines</u> (<u>DFSG</u>) part of the contract, initially designed as a set of commitments that we agree to abide by, has been adopted by the free software community as the basis of the <u>Open Source</u> <u>Definition</u>.

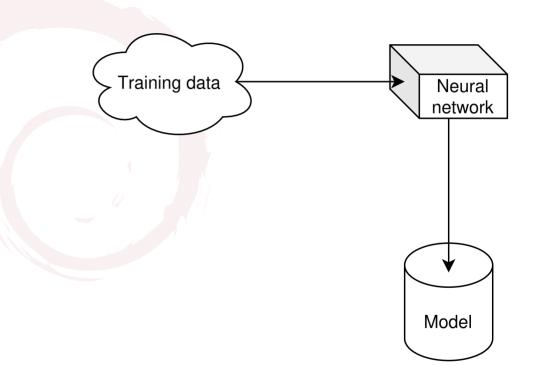
Building a software entails

- The source code
- A compiler
- The resulting program



Building an AI model entails

- Training data
- A neural network
- The resulting model



Interpretation of DFSG on Artificial Intelligence (AI) Models

Interpretation of DFSG on Artificial Intelligence (AI) Models

Al models released under DFSG-compatible license without original training data or program" are not seen as DFSG-compliant.

The "AI models released under DFSG-compatible license without original training data or program", a particular type of files as explained above, are not seen as DFSG-compliant. Hence, they can not be included in the "main" section of the Debian archive. This proposal does not specify whether the "non-free" section of Debian archive can include those files.

The point is not the definition, but the thought behind it

Interpretation of DFSG on Artificial Intelligence (AI) Models

Al models released under DFSG-compatible license without original training data or program" are not seen as DFSG-compliant.

The "AI models released under DFSG-compatible license without original training data or program", a particular type of files as explained above, are not seen as DFSG-compliant. Hence, they can not be included in the "main" section of the Debian archive. This proposal does not specify whether the "non-free" section of Debian archive can include those files.

Consequences

Requiring training data has practical, technical, and ethical consequences

Interpretation of DFSG on Artificial Intelligence (AI) Models

Al models released under DFSG-compatible license without original training data or program" are not seen as DFSG-compliant.

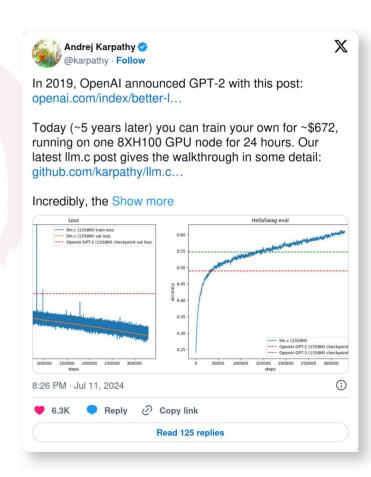
The "AI models released under DFSG-compatible license without original training data or program", a particular type of files as explained above, are not seen as DFSG-compliant. Hence, they can not be included in the "main" section of the Debian archive. This proposal does not specify whether the "non-free" section of Debian archive can include those files.

- Training dataset size
- Training time
- Reproducibility
- Security
- Transparency
- Dataset licensing

- Training dataset size
- Training time
- Reproducibility
- Security
- Transparency
- Dataset licensing

 Data is "too big", LLMs are out

- Data is "too big", LLMs are out ...for now?
- The llm.c project trained GPT-2 (124M) in 20 minutes
- The original dataset was 40 GB, the largest package is 1 GB
- TinyStories is ~2GB



- Training dataset size
- Training time
- Reproducibility
- Security
- Transparency
- (Dataset licensing)

• Secure code delivery is the problem of getting software from its author to its users safely, with a healthy dose of mistrust towards the author and everything else in between

- Secure code delivery is the problem of getting software from its author to its users safely, with a healthy dose of mistrust towards the author and everything else in between
- Build reproducibility is one vertex of the triangle





- Training dataset size
- Training time
- Reproducibility
- Security
- Transparency
- (Dataset licensing)

 In some fields, access to the training data is a crucial part of the model's value

 In some fields, access to the training data is a crucial part of the model's value I encountered an example wherein an AI model identified a previously unknown biomarker associated with cancer. This discovery was only possible because the researchers had access to the underlying dataset. Without that access, the model's findings would have been unverifiable.

_ -

Arian Ott <arian.ott@ieee.org>

 In some fields, access to the training data is a crucial part of the model's value I encountered an example wherein an AI model identified a previously unknown biomarker associated with cancer. This discovery was **only** possible because the researchers had access to the underlying dataset. Without that access, the model's findings would have been unverifiable.

_ -

Arian Ott <arian.ott@ieee.org>

Ethical consequences

- Dataset licensing (!)
 - Is copyright too strict?
- Bias of the model
- Power dynamics
 - Is it fair for the author only to know the data?

Takeaways

- LLMs are currently too large to be re-trained by Debian
 - Not all AI models are LLMs!
- Our priorities are our users and free software, not new technology
- We should focus on distributing AI models which align with our values

Why you cannot "apt install" an LLM (yet)

Thank you!

Questions?