



Guía de Buenas Prácticas para el Desarrollo y Mantenimiento de Software



GOBIERNO DE ESPAÑA

MINISTERIO DE CIENCIA E INNOVACIÓN



Instituto de Salud Carlos III

IMPACT

Infraestructura de Medicina de Precisión asociada a la Ciencia y la Tecnología

Programa	IMPACT: Infraestructura de Medicina de Precisión asociada a la Ciencia y la Tecnología		
Nombre Proyecto	IMPACT-Data: Programa de Ciencia de Datos de IMPACT		
Expediente	IMP/00019		
Duración	Enero 2021 – Diciembre 2023		
Página web	impact-data.bsc.es		
Paquete Trabajo	WP2 – INFRAESTRUCTURA COMPUTACIONAL EN LA NUBE PARA LA GESTION E INTEGRACION DE DATOS.		
Tarea	T2.2 - Adopción e implantación de protocolos para el uso de contenedores software y flujos de trabajo siguiendo las recomendaciones de ELIXIR, EOSC-Life y GA4GH		
Entregable	E2.4. Guía de Buenas Prácticas para el Desarrollo y Mantenimiento de Software		
Versión	1.1.1		
Fecha Entrega	30/06/2022	Fecha Aprobación	17/05/2023
Responsable	BSC (Organización)		
Nivel Diseminación	X	PU	Público
		CO-IMP	Confidencial, sólo participantes de los pilares de IMPACT, incluyendo la comisión de evaluación de IMPACT.
		CO-DATA	Confidencial, sólo participantes de IMPACT-Data, incluyendo la comisión de evaluación de IMPACT.

<i>Autores</i>		
<i>Organización</i>	<i>Nombre</i>	<i>Rol</i>
BSC	Lidia López	Autora
BSC	Salvador Capella-Gutierrez	Autor
CNB	Laura del Cano	Revisora
UPV	Ignacio Blanquer	Revisor

<i>Historial de versiones</i>			
<i>Nro.</i>	<i>Fecha</i>	<i>Descripción</i>	<i>Autor</i>
v 0.0	12/04/2022	Documento creado	L.López (BSC-CNS)
v 0.1	04/05/2022	primera versión índice	L.López (BSC-CNS)
v 0.2	13/05/2022	Revisada v 0.1	L. del Caño (CNB-CSIC)
v 0.3	23/05/2022	Revisada v 0.2	I. Blanquer (UPV)
v 0.4	31/05/2022	Contenido sección 2 (recomendaciones)	L.López (BSC-CNS)
v 0.5	02/06/2022	Contenido sección 3 (implementación) y anexos	L.López (BSC-CNS)
v 0.6	03/06/2022	Conclusiones	L.López (BSC-CNS)
v 0.7	10/06/2022	Comentarios sobre EDAM	L.López (BSC-CNS)
v 0.8	20/06/2022	Revisados algunos números de la sección 4	L.López (BSC-CNS)
v 0.9	24/06/2022	Revisada	L. del Caño (CNB-CSIC)
v 0.10	29/06/2022	Revisada	I. Blanquer (UPV)
v 0.11	10/07/2022	Cambios relacionados con los comentarios de los revisores	L.López (BSC-CNS)
V 0.12	14/07/2022	Revisada	Laura del Caño (CNB-CSIC)
V 0.13	21/07/2022	Lista herramientas identificadas incluida como anexo	L. López (BSC-CNS)
v 1.0	29/07/2022	Revisión general	S. Capella-Gutierrez (BSC-CNS)
v 1.1	17/05/2023	Cambio visibilidad a público y aprobado	Comité Dirección
v 1.1.1	14/06/2023	Cambio de formato para publicar en la Web de IMPaCT	David Velasco (ISCIII)

Contenido

Contenido	4
Tablas	4
Figuras	5
Resumen Ejecutivo	6
Introducción	8
Audiencia	8
Ámbito	8
Relación con otros Entregables	8
Estructura Entregable	8
1 Motivación	9
2 Buenas prácticas para el desarrollo de software de investigación de código abierto (4OSS)	9
2.1 Recomendación: Código Fuente Abierto (OSS)	10
2.2 Recomendación: Software Localizable	11
2.3 Recomendación: Escoger Licencia Apropriada	14
2.4 Recomendación: Definir Comunicación, Gobernanza y Contribución	15
3 Calidad del Software	16
4 Implementación de 4OSS en IMPaCT-Data	17
5 Conclusiones	21
Referencias	22
Anexo A. Relación de herramientas seleccionadas para Infraestructura IMPaCT-Data	23
Anexo B. Seminario bio.tools and EDAM	26
Anexo C. Seminario de encapsulado de componentes software en contenedores	32
Anexo D. Seminario Workflows - Uso de gestores de workflows, registro y compartición	38

Tablas

Tabla 1. Combinación de licencias OSS.....	15
Tabla 2. Repositorios para el catálogo de software de IMPaCT-Data.....	19
Tabla 3. Comentarios sobre la ontología EDAM para anotación de herramientas software .	20

Figuras

Figura 1. Registro bio.tools	12
Figura 2. Clasificación concepto Relation extraction en la ontología EDAM	13
Figura 3. Seminarios para aplicar 4OSS en IMPaCT-Data	18
Figura 4. Asistencia a seminarios/tutoriales	18
Figura 5. Uso de licencias OSS para las herramientas software	19
Figura 6. Uso de la ontología EDAM en El registro en bio.tools.....	20

Resumen Ejecutivo

El Software de investigación es un componente esencial en disciplinas científicas con un gran volumen de datos. En el campo de Ciencias de la Vida es difícil imaginar llevar a cabo la mayoría de las actividades científicas sin contar con el uso de software para el procesamiento, análisis, visualización e interpretación de grandes cantidades de datos. En este documento se presentan las cuatro recomendaciones incluidas en la guía de buenas prácticas para el desarrollo de software de código abierto producidas por ELIXIR. Estas recomendaciones, están basadas en los valores del código fuente abierto, y se han escrito para que el software desarrollado por investigadores sea más localizable (fácil de encontrar), reusable y transparente, en resumen, para que el software sea más fácil de mantener y extender en el tiempo. Estas cuatro recomendaciones son: (1) hacer público el código fuente en sistemas de control de versiones, p. ej. GitHub, GitLab, desde el primer día; (2) hacer el software localizable a través de su publicación en registros utilizados por la comunidad para tal fin; (3) escoger una licencia reconocida por la comunidad, con preferencia por licencias de software libre, teniendo en cuenta las licencias de los componentes del que depende el software en cuestión; y (4) definir los mecanismos de comunicación, gobernanza y colaboración más allá de los propios desarrolladores del software de investigación.

En este documento también se ha incluido las acciones que se han llevado a cabo en el marco de IMPaCT-Data para facilitar la adopción de las cuatro recomendaciones anteriores. Estas acciones se refieren a una serie de sesiones de trabajo y seminarios para facilitar acceso a las plataformas, buenas prácticas y guías a los miembros de IMPaCT-Data que les permitan implementar las recomendaciones aquí descritas. Los seminarios que se han impartido han incluido los siguientes temas: (1) registro de herramientas en el registro bio.tools y el uso de la ontología EDAM para describir el software de investigación; (2) guía de buenas prácticas para el desarrollo de software recogiendo ejemplos prácticos de cada una de las recomendaciones; (3) encapsulado de software de investigación de forma modular en contenedores; y (4) uso de gestores de flujos de trabajo, popularmente conocidos como workflows o pipelines; el registro de dichos workflows y su publicación para el potencial uso por parte de la comunidad científica. Estos seminarios han sido impartidos por ponentes nacionales e internacionales, incluyendo el Centro de Regulación Genómica y el Barcelona Supercomputing Center – Centro Nacional de Supercomputación así como la Universidad de Bergen (Noruega), la Universidad Técnica de Dinamarca (Dinamarca), la Universidad Libre de Freiburg (Alemania), todos ellos miembros de la plataforma de herramientas de ELIXIR y otras comunidades relevantes como Galaxy y Bioconda. A estos seminarios han asistido representantes de 34 de las 47 instituciones miembros del proyecto.

Como resultado de las acciones realizadas para la implantación de las recomendaciones, al cierre de este documento, se han identificado 75 herramientas dentro del consorcio, de las cuales 57 se han clasificado como herramientas software, 15 como workflows y 3 como bases de datos. De las herramientas identificadas, 51 tienen licencia OSS y 58 han sido incluidas en el registro de ELIXIR bio.tools.

Introducción

Audiencia

Este documento está destinado a todos los participantes del proyecto IMPaCT-Data, para que puedan seguir las recomendaciones para el desarrollo y mantenimiento de software de investigación. El cumplimiento de las recomendaciones es opcional dado que existe software de investigación que existe con antelación a IMPaCT-Data el cual puede utilizarse en el contexto del proyecto y cuya autoría puede ser de miembros ajenos al consorcio. En el caso de software desarrollado desde cero en el contexto de este proyecto, se espera el cumplimiento de estas recomendaciones para contribuir a su uso, difusión y sostenibilidad.

Ámbito

Las recomendaciones que se describen en este documento se aplicaran a lo largo de todo el proyecto.

Relación con otros Entregables

Este entregable no tiene relación directa con ningún otro, aunque es esperable que estas recomendaciones se sigan para el software desarrollado en el contexto de IMPaCT-Data.

Estructura Entregable

En la Sección 1 se explica la motivación para el seguimiento de las recomendaciones que forman parte de esta guía de buenas prácticas. Las recomendaciones se explican en la Sección 2. La Sección 3 incluye material complementario referente a la mejora de la calidad del código fuente de los componentes software. En la Sección 4 se han incluido las acciones realizadas hasta el momento para fomentar la implementación de las recomendaciones. Finalmente, en la Sección 5 se incluyen las conclusiones.

En los anexos de este documento se incluye la lista de herramientas identificadas para el catálogo de IMPaCT-Data (Anexo A) y el material de soporte para implementar algunas de las recomendaciones. El Anexo B incluye detalles sobre el registro de herramientas software en el registro de ELIXIR bio.tools. El Anexo C incluye detalles sobre el encapsulado de software en contenedores para la instalación de herramientas software en distintos entornos computacionales. Finalmente, el Anexo D incluye detalles sobre la gestión de flujos de trabajo (*workflows* o *pipelines*), en concreto sobre distintos gestores de workflows, su registro y publicación para el potencial uso por parte de la comunidad, ya sea de IMPaCT-Data o ajena al proyecto.

1 Motivación

El Software es una parte importante en la investigación en el campo de ciencias de la vida. De hecho, ciertas investigaciones no se podrían realizar si no existieran las herramientas software que lo permitieran. Cuando el software se desarrolla como parte de las herramientas que dan soporte a la investigación, habitualmente no se considera como un resultado de la investigación, por lo que no se le aplican los mismos procesos para asegurar su calidad, referencia, reproducibilidad y reusabilidad.

La producción y uso de software de código abierto (OSS, siglas en inglés para Open Source Software) es uno de los mecanismos que permiten aplicar los principios de ciencia abierta a los componentes software. OSS es software para el cual el código fuente es públicamente accesible para cualquier persona que quiera inspeccionar, usar, modificar y mejorar dicho código. Tener el código de las herramientas software que se utilizan para la investigación contribuye al reconocimiento de los desarrolladores del mismo, a crear comunidad alrededor de dichos desarrollos, así como a aumentar la confianza sobre el mismo al permitir la posibilidad de inspeccionar dicho código. Es importante señalar que, dependiendo de las licencias de los componentes utilizados en el desarrollo de software, puede ser necesario utilizar el mismo tipo de licencias, p. ej. familia de licencias GPL.

En este entregable se exponen recomendaciones para el desarrollo de software de código abierto en el contexto de ciencia abierta.

2 Buenas prácticas para el desarrollo de software de investigación de código abierto (4OSS)

Estas recomendaciones están propuestas en la guía de buenas prácticas definida por el grupo de trabajo *Software development best practices for Life Sciences*¹ de la organización ELIXIR² [1]. El objetivo de este grupo de trabajo es mejorar la calidad y sostenibilidad del software desarrollado por investigadores en el dominio de ciencias de la vida.

Estas recomendaciones se resumen en:

- Hacer público el código fuente en sistemas de control de versiones, p. ej. GitHub, GitLab, desde el primer día.
- Hacer el software localizable a través de su inclusión en registros generales o específicos de la comunidad en cuestión.
- Escoger la licencia más adecuada.
- Definir Comunicación, Gobernanza y Contribución.

¹ <https://elixir-europe.org/platforms/tools/software-best-practices>

² <https://elixir-europe.org/>

Estas recomendaciones no pretenden substituir las guías de buenas prácticas existentes para el desarrollo de software, sino complementarlas. Estas recomendaciones están basadas en los valores del código fuente abierto, y se han escrito para mejorar el software desarrollado por investigadores permitiendo que sea más localizable (fácil de encontrar), re-usable y transparente.

De hecho, estas recomendaciones están alineadas con los principios FAIR (siglas en inglés para los términos **F**indable, **A**ccesible, **I**nteroperable y **R**e-usable) definidos en el contexto de gestión de datos en el ámbito científico, y en la actualidad trasladados a distintos tipos de objetos digitales, incluido el software de investigación.

Además de la definición y publicación de estas recomendaciones, el grupo de trabajo *Software development best practices for Life Science*, en colaboración con *the Carpentries*³, ha producido el material de formación *4 Simple recommendations for Open Source Software (4OSS)*⁴.

2.1 Recomendación: Código Fuente Abierto (OSS)

Para maximizar la reproducibilidad, reusabilidad y la colaboración se recomienda tener el código fuente en un repositorio abierto (de acceso público) con control de versiones, por ejemplo, GitHub o GitLab. Tener el código fuente disponible públicamente también ayuda a mejorar la visibilidad y la confianza en el componente software.

Los beneficios de esta recomendación son:

- Promueve la confianza del componente software en concreto y del proyecto en el que se está desarrollando en general.
- Facilita encontrar proyectos existentes en los que se está desarrollando software.
- El uso de un sistema de control de versiones proporciona acceso al historial de las contribuciones, lo cual ayuda a dar reconocimiento a los contribuidores de código.
- Incentiva la contribución de la comunidad.
- Incrementa las oportunidades de colaboración y re-uso.
- Expone el trabajo a la evaluación de la comunidad, posibilitando sugerencias y validaciones del código.
- Incrementa la transparencia para el escrutinio de la comunidad.
- Incentiva a los desarrolladores a pensar e implementar buenas prácticas al desarrollar código.
- Facilita la reproducción de los resultados científicos generados por las versiones anteriores del software.
- Incentiva a los desarrolladores a generar documentación, incluyendo un manual de usuario y comentarios en el código fuente.

³ <https://carpentries.org/>

⁴ <https://softdev4research.github.io/4OSS-lesson/>

- Garantiza el almacenamiento a largo plazo del código fuente.

Una buena práctica al crear un repositorio de código en un repositorio abierto es incluir el archivo *README.md*, en el directorio raíz, incluyendo una descripción del componente para que los usuarios o potenciales contribuidores puedan entender su propósito y/o funcionalidad, así como las políticas de contribución al código.

2.2 Recomendación: Software Localizable

Consiste en registrar el software en un registro que sea popular en el dominio, al registrarlo incluir los suficientes metadatos para maximizar las posibilidades que sea encontrado en las búsquedas.

Los beneficios de esta recomendación son:

- Incrementa la visibilidad del proyecto, del software, su uso, sus éxitos, referencias y quienes contribuyen.
- Incentiva a los desarrolladores de software a pensar que metadatos describen mejor el software y cómo exponerlo.
- Ayuda a exponer los metadatos en un formato entendible por las máquinas (p. ej. Buscadores) a través del registro.
- Incrementa las opciones de colaboración, re-uso y mejora.

Dependiendo del tipo de software, hay que buscar el registro adecuado para darle la máxima visibilidad en el dominio más adecuado. Por ejemplo, en el contexto de IMPaCT-Data, para las herramientas software hemos escogido el registro bio.tools y para los flujos de trabajo, el registro WorkflowHub. Estos dos recursos forman parte del ecosistema de ELIXIR, lo cual se alinea con las directivas fijadas por IMPaCT. Cada registro tiene sus normas para la anotación de sus entradas, por ejemplo en el registro bio.tools se utiliza la ontología EDAM⁵ para la anotación de los recursos que contiene (ver Anexo B), que contiene más de 3.500 conceptos clasificados y relacionados entre ellos. De hecho, en EDAM los conceptos están clasificados como:

- Tema (traducción de *Topic*⁶): Una categoría que denota un dominio o campo de interés, de estudio, aplicación, trabajo, datos o tecnología. Los temas no tienen fronteras claramente definidas entre sí.
- Operación (*Operation*⁷): Una función que procesa un conjunto de entradas y da como resultado un conjunto de salidas, o asocia los argumentos (entradas) con los valores

⁵ <https://edamontology.org>

⁶ https://bioportal.bioontology.org/ontologies/EDAM?p=classes&conceptid=topic_0003

⁷

https://bioportal.bioontology.org/ontologies/EDAM/?p=classes&conceptid=http%3A%2F%2Fedamontology.org%2Foperation_0004

(salidas). Datos (*Data*⁸): Información, representada en un artefacto de información (registro de datos) que es “comprensible” por herramientas computacionales dedicadas que pueden usar los datos como entrada o producirlos como salida.

- Formato (*Format*⁹): Una forma definida o diseño de representar y estructurar datos en un archivo de ordenador, blob, string, mensaje o en cualquier otro lugar.

Las anotaciones en el registro bio.tools utilizan los conceptos asociados a las categorías *Topic* y *Operation*. En la Figura 1 se muestran dos de las herramientas de IMPaCT-Data registradas en bio.tools, después de la descripción se ven los conceptos relacionados con el *Topic* en la primera línea y con el *Operation* en la segunda, las etiquetas de la tercera línea son otras clasificaciones fuera de EDAM (tipo herramienta, licencia, colección, etc.).

The screenshot shows the bio.tools website interface. At the top, there is a search bar and navigation links. The main content area displays two tool entries. The first entry is for 'IMPACT-Data', which includes a description of its purpose and a list of associated concepts. The second entry is for 'RD-Connect Genome-Phenome Analysis Platform (GPAP)', also with a description and a list of associated concepts. The concepts are displayed as colored tags below each tool's description.

Figura 1. Registro bio.tools

8

https://bioportal.bioontology.org/ontologies/EDAM/?p=classes&conceptid=http%3A%2F%2Fedamontology.org%2Fdata_0006

12

9

https://bioportal.bioontology.org/ontologies/EDAM/?p=classes&conceptid=http%3A%2F%2Fedamontology.org%2Fformat_1915

Para la herramienta *RD-Connect Genome-Phenome Analysis Platform (GPAP)* de la Figura 1, los conceptos de EDAM son:

- Topic: *Genotype and phenotype* (de la clasificación *Biology/Genetics*), *Rare diseases (Medicine/Pathology)*, *Medical Informatics (Informatics classification)*, y *Data Mining (Computer Science)*
- Operation: *Data Retrieval* (clasificaciones¹⁰ *Analysis/Text mining/Information Retrieval* and *Data Handling/Query and Retrieval*), *Information Extraction (Analysis/Text mining* y *Prediction and recognition/Text mining)*, y *Relation Extraction (Analysis/Text mining* y *Prediction and recognition/Text mining)*.

En la Figura 2 se puede ver la clasificación del concepto *Relation extraction* que pertenece a dos clasificaciones dentro de *Operation*.

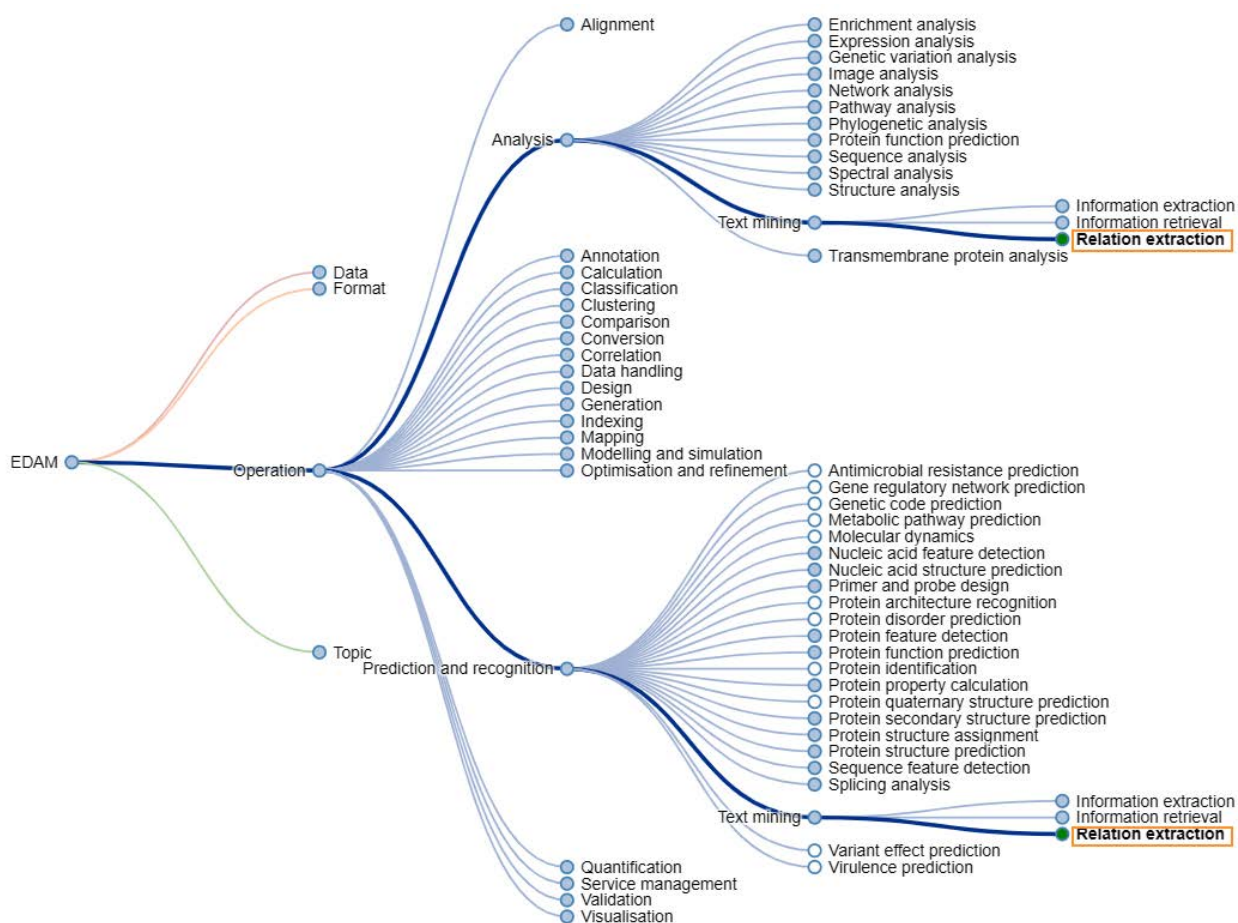


Figura 2. Clasificación concepto Relation extraction en la ontología EDAM

¹⁰ Un concepto puede estar clasificado en varias clasificaciones

2.3 Recomendación: Escoger Licencia Apropiaada

A la hora de desarrollar software de investigación, es importante asignar una licencia a dicho software. Siguiendo las directrices marcadas por IMPaCT, el software desarrollado debe estar disponible bajo licencias de código abierto (Open Source). Escoger la licencia Open Source adecuada permite clarificar los términos y condiciones de uso, modificación y redistribución del código fuente, así como fija responsabilidades de los autores del software. En el caso que el software utilice otros componentes disponibles bajo licencias de código abierto, se debe comprobar que se están cumpliendo los términos y condiciones de la licencia de cada componente y si la licencia es compatible con la que hemos escogido para el nuestro.

Los beneficios de esta recomendación son:

- Clarifica las responsabilidades y derecho de terceros que quieran utilizar, copiar, redistribuir, modificar o re-usar el código fuente.
- Permite utilizar el código fuente en jurisdicciones en las que “código sin licencia” significa que no se puede utilizar de ninguna manera.
- Protege la propiedad intelectual del software.
- Proporciona sostenibilidad a largo plazo permitiendo contribuciones y reutilización en el contexto de organizaciones legales con financiación.

Una buena práctica al crear un repositorio en un repositorio abierto es incluir el archivo *LICENSE.md*, en el directorio raíz, incluyendo el texto de la licencia escogida para el componente.

En cuanto al uso de licencias de código abierto, se recomienda escoger una licencia aprobada por la Open Source Initiative¹¹ (OSI), excepto que se necesiten condiciones especiales. En la web de OSI se pueden encontrar la lista de las licencias aprobadas y los criterios a cumplir por las licencias.

Las licencias OSS se distinguen entre *copyleft* y permisivas. Las licencias *copyleft* se caracterizan porque los trabajos derivados deben ser puestos a disposición de la comunidad obligatoriamente, lo que significa que el software se debe redistribuir utilizando la misma licencia o una más restrictiva. Las licencias permisivas no incluyen restricciones en los trabajos derivados, siempre que se de crédito a los autores originales, lo cual permite su potencial uso comercial. Hay un tercer grupo de licencias que está entre los dos tipos anteriores, el código en sí está disponible bajo una licencia *copyleft* pero los componentes (normalmente librerías) se pueden combinar con otro tipo de licencias sin la obligación de distribución bajo el mismo tipo de licencia (las llamaremos *copyleft* permisivas).

Si queremos distribuir un componente OSS, la licencia que se puede escoger para la distribución depende de dos factores: (1) la licencia original del componente y (2) el tipo de distribución que se quiere hacer. En la siguiente tabla se definen los tipos de licencia que se pueden aplicar (celdas) a la distribución de un software dependiendo de la licencia original

¹¹ <https://opensource.org/licenses>

(columnas) y de su tipo de distribución (filas). Se ha considerado la distribución como trabajo derivado (se incluyen modificaciones en el código fuente), si se distribuye combinado con otros componentes y ambos.

Tabla 1. Combinación de licencias OSS

	Copyleft	Copyleft permisiva	Permisiva
Derivado	Sólo si la licencia es copyleft	Dependiendo de cómo se enlace (linking) el componente dentro del proyecto	Cualquier tipo de licencia
Combinado	Solo si la licencia es copyleft	Cualquier tipo de licencia	Cualquier tipo de licencia
Derivado y combinado	Solo si la licencia es copyleft	Dependiendo de cómo se enlace (linking) el componente dentro del proyecto	Cualquier tipo de licencia

2.4 Recomendación: Definir Comunicación, Gobernanza y Contribución

Definir los procesos de comunicación, gobernanza y contribución no significa que el software deba ser desarrollado colaborativamente. Sin embargo, hay que definir de manera clara la estrategia de colaboración y contribución teniendo un modelo transparente de gobernanza y canales de comunicación.

Los beneficios de esta recomendación son:

- Incrementa la transparencia en cómo el proyecto y el desarrollo del software se gestiona.
- Ayuda a definir responsabilidades y procesos de toma de decisiones.
- Ayuda a la comunidad a colaborar, contribuir y comunicarse con el proyecto.

La comunicación se refiere a establecer los mecanismos necesarios para el trabajo colaborativo dentro del proyecto, incluyendo la puesta a punto de distintos canales de comunicación. De hecho, los canales de comunicación dependerán de diversos factores teniendo en la cantidad de peticiones de mejora esperadas y del tamaño del equipo de desarrollo. Estos mecanismos se deben implantar desde el inicio del proyecto. Por ejemplo, el uso de listas de distribución, plataformas de mensajería como Slack¹² o rocket chat¹³ y los *issues* en los propios repositorios de código. Para comunicaciones importantes, como por ejemplo la publicación de una nueva versión, estas se pueden hacer a través de un blog asociado al proyecto u otro espacio web y a través de los canales de los contribuidores del código.

¹² <https://slack.com/intl/es-es/>

¹³ <https://www.rocket.chat/>

La gobernanza se refiere a los procesos de toma de decisiones, idealmente se debe definir las diferentes responsabilidades (roles) y reconocimiento. Es importante no dejar responsabilidades sin definir dado que esto puede crear incertidumbre a medida que el proyecto software avanza y crece.

Para incentivar la colaboración de agentes externos, la inclusión de las guías de contribución es muy útil, por ejemplo, añadir el archivo *CONTRIBUTING.md* en el directorio raíz repositorio. Esta guía debe incluir como mínimo:

- Reconocimiento/bienvenida a las personas que pretenden contribuir.
- Descripción de las diferentes maneras de contribución (secciones “How-to”). Por ejemplo, como reportar un error (*bug*), como contribuir con código, como sugerir mejoras. En el caso de permitir la contribución con código, habría que incluir un enlace a las convenciones de codificación y la guía de estilo.
- La manera de contactar.

En el curso producido por el grupo de trabajo de ELIXIR *Software development best practices for Life Sciences*¹⁴ (sección 4) se incluye una lista de verificación que incluye aspectos que se deben tener en cuenta para hacer las contribuciones más fáciles, claras y transparentes.

También es importante definir un código de conducta, que se puede incluir en el repositorio raíz como un archivo separado (*CODE_OF_CONDUCT.md*) o en una sección del archivo *CONTRIBUTING.md*. Las comunidades Open Source no solo son espacios en los que se comparte tecnología, tienen un componente social de comunidad. El código de conducta define lo que se espera de los miembros de la comunidad en sus interacciones, lo que ayuda a mantener el bienestar de la comunidad. En el código de conducta se espera encontrar:

- Comportamiento esperado
- Comportamiento no aceptado
- Consecuencias en caso de comportamiento no aceptado
- Protocolo para reportar un comportamiento no aceptado por parte de alguno de los miembros

3 Calidad del Software

Las recomendaciones incluidas en la sección anterior no se refieren a la calidad del código fuente, son guías para mejorar la sostenibilidad del software desarrollado en entornos de investigación. Para facilitar el re-uso y el mantenimiento de estos componentes software, el código fuente debe alcanzar niveles de calidad aceptables.

Las actividades relacionadas con la gestión del proyecto también pueden influir en la calidad del software resultante. En el grupo de trabajo *Software development best practices for Life Sciences* en ELIXIR se está trabajando en un *Software Management Plan*¹⁵ que puede ser

¹⁴ <https://softdev4research.github.io/4OSS-lesson/>

¹⁵ <https://elixir-europe.org/sites/default/files/documents/software-management-plan.pdf>

utilizado como una lista de comprobación (checklist) para comprobar si se están haciendo las actividades recomendadas para maximizar la calidad del software en desarrollo. Las actividades incluidas en este plan se refieren a documentación; pruebas (*testing*); interoperabilidad; comunidad, contribución y gobernanza, reproducibilidad y reconocimiento.

En el contexto del European Open Science Cloud¹⁶ (EOSC), el proyecto EOSC-Synergy¹⁷ tiene como objetivo el proporcionar servicios con calidad certificada. EOSC-Synergy ha definido un conjunto de criterios que se pueden utilizar como marco de referencia para maximizar la calidad de los desarrollos de software en los proyectos de investigación [2]. Estos criterios se refieren a la accesibilidad del código, selección de licencia, estilo del código, meta-datos del código, pruebas unitarias, herramientas para la automatización de las pruebas (*test harness*), desarrollo dirigido por pruebas (TDD por sus siglas en inglés *Test-Driven Development*), documentación, seguridad, modelos de ramas en el código (*code workflow/branching model*) de tal forma que correcciones y nuevas funcionalidades se desarrollan y prueban en ramas independientes, y cuando se ha alcanzado el nivel de calidad deseado se integran en la rama principal de desarrollo, versiones semánticas, gestión del código, revisión de código, entrega automática y despliegue automático.

4 Implementación de 4OSS en IMPaCT-Data

En el contexto del proyecto, se han organizado cuatro seminarios para facilitar la adopción de las recomendaciones descritas en este documento (ver Tabla 3):

- bio.tools and EDAM: How to publish information about your tools. Matúš Kalaš (Universitet of Bergen, Noruega) and Hans Ienasescu (Danmarks Tekniske Universitet, Dinamarca).
- Best practices for software development. Salvador Capella-Gutierrez (Barcelona Supercomputing Center - Centro Nacional de Supercomputación, España).
- Containarización componentes software. Björn Grüning (Bioconda/BioContainers community, Germany).
- Workflows: Uso de gestores de workflows, registro y compartición. José Espinosa-Carrasco, Center for Genomic Regulation (Centre for Genomic Regulation), Ignacio Eguinoa (Ghent University, ELIXIR-BE, Bélgica), Salvador Capella-Gutierrez, Barcelona Supercomputing Center – Centro Nacional de Supercomputación (BSC-CNS).

¹⁶ <https://eosc-portal.eu/>

¹⁷ <https://www.eosc-synergy.eu/>



Figura 3. Seminarios para aplicar 4OSS en IMPaCT-Data

De las 47 instituciones que forman parte de IMPaCT-Data, han asistido representantes de 34 instituciones (72%) a al menos una de las sesiones. En la Figura 4 se indica el número de personas que asistieron a cada sesión (serie azul y con trama de rayas) y el número de instituciones identificadas (serie naranja con trama de cuadros). Aunque todas las sesiones han tenido un número elevado de asistentes, las que han suscitado más interés han sido la de buenas prácticas y encapsulado de software en contenedores para su instalación, ambas con una asistencia de 48 personas.

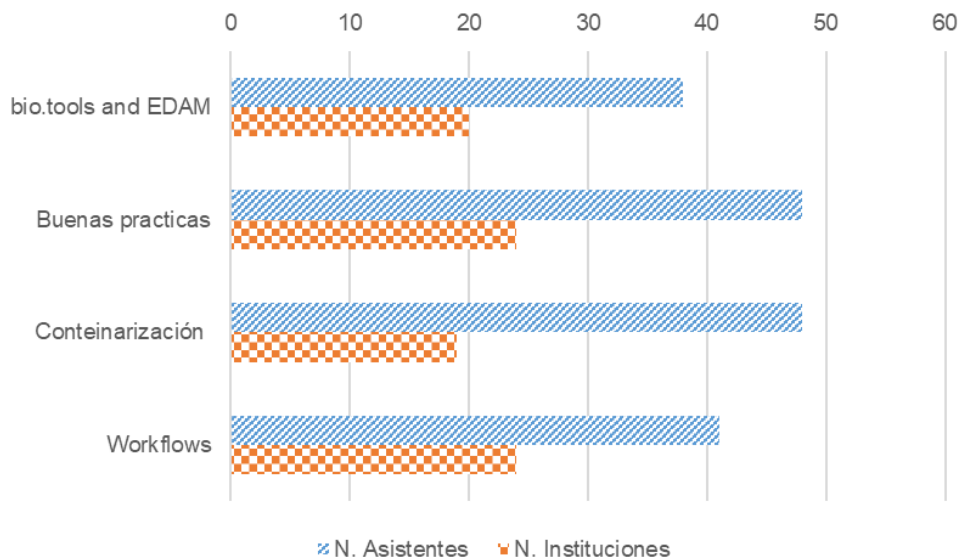


Figura 4. Asistencia a seminarios/tutoriales

Se ha pedido a los miembros del proyecto que nos indiquen las herramientas software/workflows que utilizan en el proyecto. Hasta el momento se han identificado 75 herramientas de 19 instituciones, 57 se han clasificado como herramientas software, 15 como workflows y 3 como bases de datos (listadas en el Anexo A).

Referente a la recomendación de tener el código fuente abierto (Sección 2.1), se han identificado las herramientas que se están desarrollando bajo una licencia OSS. De las 75 herramientas identificadas, se ha recogido esta información para 66 (9 herramientas no han proporcionado información sobre si la licencia es OSS o no), con el resultado de que la mayoría tienen licencia OSS (77%).

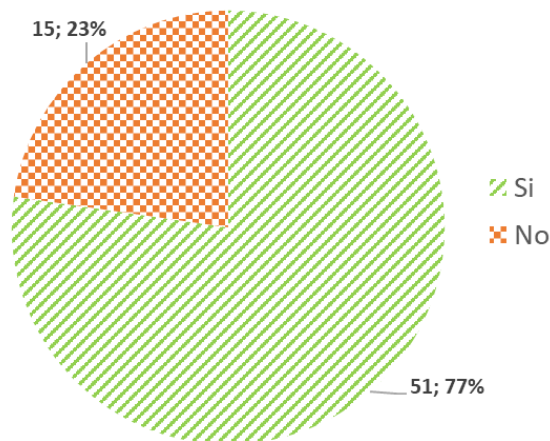


Figura 5. Uso de licencias OSS para las herramientas software

Para seguir la recomendación de hacer el software localizable (Sección 2.2, se han seleccionado diferentes repositorios que se utilizarán en el proyecto para mantener el catálogo de software. Las diferentes herramientas software pueden ser de diferente tipo y ser accesibles de diferente manera, la Tabla 2 incluye los detalles de los repositorios seleccionados.

Tabla 2. Repositorios para el catálogo de software de IMPaCT-Data

	Repositorio	Metadata	Detalles
Software	bio.tools ¹⁸	EDAM (Anexo A)	IMPACT-Data collection y domain
Containers	BIOCONDA ¹⁹	--	--
Workflows	WorkflowHub ²⁰	EDAM (Topic, Operation)	--

De las 75 herramientas identificadas, a la finalización de este entregable, 58 se han registrado en el registro bio.tools. De las 58 registradas, 54 se han incluido en el dominio IMPaCT-Data²¹, las 4 herramientas que no se han podido incluir en el dominio es debido a algún problema técnico que se está solucionando.

Se ha consultado a los miembros del proyecto por el uso que han hecho de la ontología EDAM para añadir los metadatos a sus herramientas al registrarlas, en concreto, se les ha preguntado si los términos definidos eran suficientes. De las 58 herramientas registradas, se

¹⁸ <https://bio.tools/>

¹⁹ <https://bioconda.github.io/>

²⁰ <https://workflowhub.eu/>

²¹ <https://bio.tools/t?domain=impact-data>

han recibido 55 respuestas, en las que para la mayoría (73%) los términos definidos en la ontología han sido suficientes.

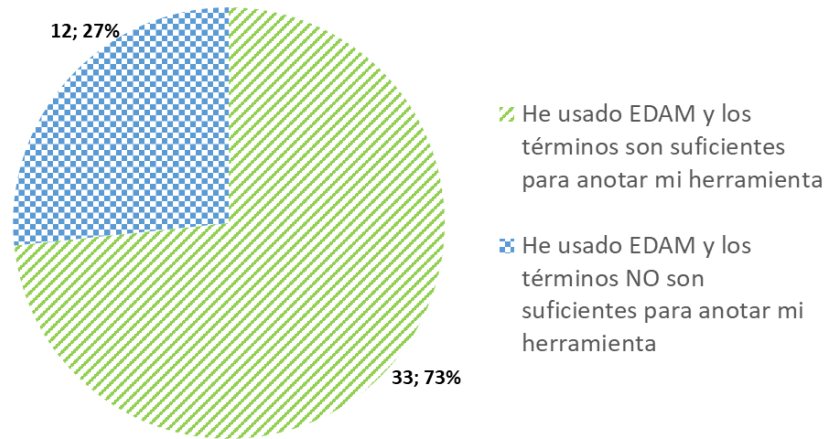


Figura 6. Uso de la ontología EDAM en El registro en bio.tools

Se han identificado a las personas de contacto de las herramientas que han reportado insuficiencias en la ontología EDAM, los comentarios recibidos se incluyen en la Tabla 3.

Tabla 3. Comentarios sobre la ontología EDAM para anotación de herramientas software

Herramienta	URL en bio.tools	Comentario sobre EDAM
mCSEA	https://bio.tools/mcsea	Han faltado términos sobre datos de metilación como input
ImaGEO	https://bio.tools/ImaGEO	Han faltado términos sobre meta-análisis
MetaGenyo	https://bio.tools/metagenyo	Han faltado términos sobre meta-análisis
DatAC	https://bio.tools/datac	Han faltado términos sobre epidemiología y factores ambientales
DExMA	https://bio.tools/dexma	Han faltado términos sobre meta-análisis
DomFun	https://bio.tools/domfun	Han faltado términos relacionados con transcriptómica (RNA-seq, miRNA-seq)
ExpHunterSuite	https://bio.tools/exphuntersuite	Han faltado términos relacionados con transcriptómica (RNA-seq, miRNA-seq)
iSkyLIMS	https://bio.tools/iskylims	Han faltado términos relacionados con la gestión de datos de laboratorio, o plataforma web de gestión de datos
Taranis	https://bio.tools/taranis	Han faltado términos relacionados con wg/cgMLST, MLST

5 Conclusiones

En este documento se presentan las cuatro recomendaciones incluidas en la guía de buenas prácticas para el desarrollo de software de código abierto producidas por ELIXIR. Estas recomendaciones están basadas en los valores del código fuente abierto, se han escrito para que el software desarrollado por investigadores sea más localizable (fácil de encontrar), reusable y transparente. Estas cuatro recomendaciones son: (1) Tener el código fuente abierto desde el primer día, (2) Hacer el software localizable; (3) Escoger la licencia más adecuada y (4) definir comunicación, gobernanza y colaboración.

Estas recomendaciones se han complementado con herramientas que se pueden utilizar para mejorar la calidad el código fuente.

En este documento también se ha incluido las acciones que se han llevado a cabo en el proyecto para facilitar la adopción de las cuatro recomendaciones. Estas acciones se refieren a una serie de sesiones de trabajo y seminarios para dar las herramientas a los miembros del proyecto para que las puedan implementar.

Como resultado de estas sesiones, al cierre de este documento, se han identificado 75 herramientas, de las cuales 57 se han clasificado como herramientas software, 15 como workflows y 3 como bases de datos. De las herramientas identificadas, 51 tienen licencia OSS y 58 han sido registradas en el portal bio.tools.

Referencias

[1] Jiménez, R. C., Kuzak, M., Alhamdoosh, M., Barker, M., Batut, B., Borg, M., ... & Crouch, S. (2017). Four simple recommendations to encourage best practices in research software. *F1000Research*, 6.

[2] Orviz, P., Lopez, A., Duma, D. C., David, M., Gomes, J., & Donvito, G. (2021). *A set of Common Software Quality Assurance Baseline Criteria for Research Projects*. Manubot.

Anexo A. Relación de herramientas seleccionadas para Infraestructura IMPaCT-Data

En este anexo se incluye la lista de 75 herramientas identificada en el contexto de IMPaCT-Data, de esta lista 58 están registradas en el registro bio.tools siguiendo las recomendaciones incluidas en este documento (ver sección 2.2)

Nombre	Tipo	OSS	URL bio.tools	Institución
ADEx	Software	Si	https://bio.tools/ADEx	FPS
APID Interactomes	Software	Si	https://bio.tools/apid	CSIC
APPRIS	Software		https://bio.tools/appris	CNIO
Automatic segmentation tool	Software	Si	https://bio.tools/automatic_segmentation_tool	INIBICA
Beyondcell	Software		https://bio.tools/beyondcell	CNIO
CohortAnalyzer	Workflow	Si		UMA
CoV-hipathia	Software	Si	https://bio.tools/cov-hipathia	FPS
CSVS	Software	Si	https://bio.tools/csvs	FPS
cypathia	Software	Si	https://bio.tools/cypathia	FPS
DatAC	Software	Si	https://bio.tools/datac	FPS
DExMA	Software	Si	https://bio.tools/dexma	FPS
DisGeNET	Software	Si	https://bio.tools/disgenet	IMIM
DiSMed	Software	Si	https://bio.tools/dismed	FISABIO
DomFun	Software	Si	https://bio.tools/domfun	UMA
DREIMT	Software		https://bio.tools/dreimt	CNIO
EvolClust	Database	No	https://bio.tools/EvolClust	BSC
ExpHunterSuite	Software	Si	https://bio.tools/exphuntersuite	UMA
FAIR4Health Data Curation Tool	Software	Si	https://bio.tools/fair4health_data_curation_tool	SAS-HUVR
FAIR4Health Data Privacy Tool	Software	Si	https://bio.tools/fair4health_data_privacy_tool	SAS-HUVR
FireDB	Software		https://bio.tools/firedb	CNIO
FJD-pipeline	Software	Si	https://bio.tools/fjd-pipeline	IIS-FJD
GeneCodis	Software	No	https://bio.tools/genecodis	FPS
GLOWgenes	Software	Si	https://bio.tools/glowgenes	IIS-FJD
GPAP	Software	No	https://bio.tools/rd-connect_platform	CRG
hipathia	Software	Si	https://bio.tools/hipathia-gemomics	FPS
Hipathia	Software	Si	https://bio.tools/hipathia	FPS
ImaGEO	Software	Si	https://bio.tools/ImaGEO	FPS
ImpuSARS	Software	Si	https://bio.tools/impusars	FPS
IonGAP	Software	Si	https://bio.tools/iongap	FIISC

iSkyLIMS	Software	Si	https://bio.tools/iskylims	ISCIII
Jupyter Hub	Software			FPS

Relación de herramientas seleccionadas para Infraestructura IMPaCT-Data (continuación)

<i>Nombre</i>	<i>Tipo</i>	<i>OSS</i>	<i>URL bio.tools</i>	<i>Institución</i>
Liferay	Software	Si	https://bio.tools/liferay	HCB
LinKEHR	Software	Si	https://bio.tools/linkehr	HCB
mCSEA	Software	Si	https://bio.tools/mcsea	FPS
Metabolizer	Software	Si	https://bio.tools/metabolizer	FPS
MetaFun	Software	Si	https://bio.tools/metafun	CIPF
MetaGenyo	Software	No	https://bio.tools/metagenyo	FPS
meTAline	Workflow	Si		BSC
MetaPhors	Database	No	https://bio.tools/metaphors	BSC
MIDS	Software	Si		FISABIO
MIGNON	Software	Si	https://bio.tools/mignon	FPS
Mini-IsoQLR	Software	Si	https://bio.tools/mini-isoqlr	IIS-FJD
MyPROSLE	Software	Si	https://bio.tools/myprosele	FPS
NanoCLUST	Software	Si	https://bio.tools/nanoclust	FIISC
NanoDJ	Software	Si	https://bio.tools/NanoDJ	FIISC
NanoRtax	Software	Si	https://bio.tools/nanortax	FIISC
NetAnalyzer	Software	Si		UMA
nf-core-viralrecon	Workflow	Si	https://bio.tools/nf-core-viralrecon	ISCIII
ngsCAT	Software	Si	https://bio.tools/ngscat	FPS
OpenEBench	Software	Si	https://bio.tools/openebench	BSC
PanDrugs	Software		https://bio.tools/pandrug	CNIO
PhenCo	Workflow	Si		UMA
PhenFun	Workflow	Si		UMA
PhylomeDB	Database	No	https://bio.tools/PhylomeDB	BSC
Pipeline CNVs germinal	Workflow	No		IdiPaz
Pipeline Germinal v1	Workflow	No		IdiPaz
Pipeline Germinal v2	Workflow	No		IdiPaz
Pipeline identificación de fusiones génicas en RNASeq	Workflow	No		INCLIVA
Pipeline Mosaico	Workflow	No		IdiPaz

Pipeline paneles/exomas en muestras tumorales	Workflow	No		INCLIVA
Pipeline RNASeq	Workflow	No		INCLIVA
Pipeline RNA-Seq	Workflow	No		IdiPaz
Pipeline Somáticas v1	Workflow	No		IdiPaz
PlasmidID	Software	Si	https://bio.tools/plasmidid	ISCI

Relación de herramientas seleccionadas para Infraestructura IMPaCT-Data (continuación)

<i>Nombre</i>	<i>Tipo</i>	<i>OSS</i>	<i>URL bio.tools</i>	<i>Institución</i>
PriorR	Software	Si	https://bio.tools/priorr	IIS-FJD
PTMCode	Software	Si	https://bio.tools/ptmcode	IIS-FJD
RStudio Workbench	Software			FPS
Slides Viewer	Software	Si	https://bio.tools/slides_viewer	INIBICA
SMAca	Software	Si	https://bio.tools/smaca	FPS
SPACNACS	Software	Si	https://bio.tools/spacnacs	FPS
Taranis	Software	Si	https://bio.tools/taranis	ISCI
TFTEA	Software	Si	https://bio.tools/tfta	FPS
TRIFID	Software		https://bio.tools/trifid	CNIO
vulcanSpot	Software		https://bio.tools/vulcanSpot	CNIO
XICRA	Workflow	Si	https://bio.tools/xicra	IGTP

Anexo B. Seminario bio.tools and EDAM



bio.tools and EDAM:
How to publish information about your tools

Hans Ienasescu and Matúš Kalaš, 17 February 2022

www.elixir-europe.org

The ELIXIR ToolsPlatform

Goal: Improve the **findability, quality, and sustainability** of software tools.

- Helps life scientists find, deploy, and compare tools, including workflows.
- Helps software providers and developers develop better software tools, describe them, and integrate them into workflows.



EDAM ontology

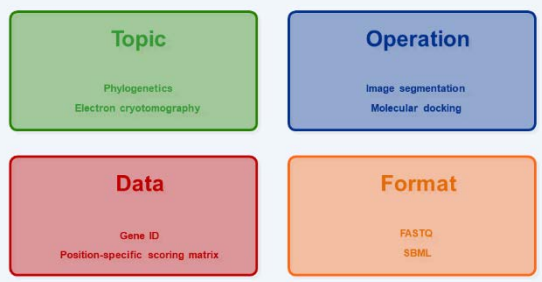
What is EDAM?

~3500 concepts in data analysis and management ...

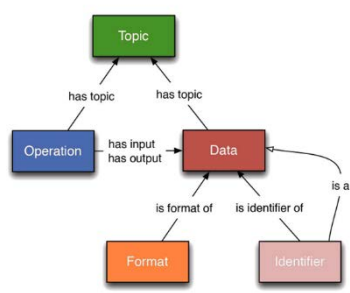
... with definitions, relations, synonyms, etc.

- Enrichment analysis
- Expression analysis
- Genetic variation analysis
- Image analysis
- Pathway or network analysis
- Phylogenetic tree analysis
- Protein function prediction
- Sequence analysis
- Spectral analysis
- Structure analysis
- Text mining
- Information extraction
- Information retrieval
- transmembrane protein analysis
- Annotation
- Calculation
- Isotopic distributions calculation
- Nucleic acid property calculation
- Protein property calculation
- Rarefaction
- Retention time prediction
- Sequence composition calculation
- Clustering
- URI
- Image
- Keyword
- Map
- Map data
- Mathematical model
- Matrix
- Molecular property
- Molecular simulation data
- Ontology data
- Over-representation data
- Pathway or network
- Phylogenetic data
- Reaction data
- Regular expression
- Report
- Scores
- Sequence
- Sequence attribute
- Sequence coordinates
- Sequence features
- Sequence features metadata
- Sequence signature data
- Sequence variations
- Binary format
- INIS
- zbit
- ABI
- ABI
- ABI
- BAI
- BAM
- BCT
- big5
- big5ed
- big5g
- Phylo
- BMP
- BTrack
- COMBINE OMEG
- OSAM
- DICOM format
- ebwt
- ebwt1
- GF
- HDF
- HDF5
- lib
- IDAM
- im

Scope of EDAM, and example concepts



Relations between concepts in EDAM



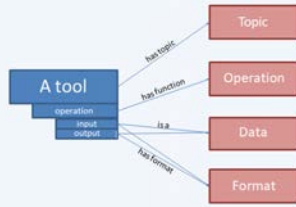
Usage areas of EDAM

- Searching for tools, workflows, learning materials, ...
- Data provenance (metadata)
- Tools and data integration
- Text mining
- Choosing terminology

EDAM became a ubiquitous component of numerous resources

... and many more ...

Annotation of computational tools with EDAM



bio.tools

The registry's motivation

- Large number of tools and databases created to support (life science) research
- Fragmented documentation and access of resources:
 - Hard to find, understand, compare and use
 - Lack formalized descriptions of scientific and technical functions
 - Don't ensure persistence of identifiers
 - Make reproducibility hard
- Need for a registry/portal/database of life science software tools
- Examples of tools registries/portals
 - EMBOSS
 - EMBRACE
 - Bio-Catalogue
 - SEQanswers
 - Omics Tools
 - BioMedBridges
 - Debian Med
 - Bioconductor

What is bio.tools?

- bio.tools strives to provide a comprehensive registry of software and data services facilitating researchers from across the spectrum of biological and biomedical science to find, understand, utilise and cite the resources they need in their day-to-day work

Created in the context of the [ELIXIR Europe](#) life science infrastructure project

<https://bio.tools>

From simple command-line tools and online services, through to databases and complex, multi-functional analysis workflows

Contains tool descriptions (information, annotations about tools), not the actual tools themselves

Open Data	Open Source	Built by You
Open data means that the data is freely available to all under CC BY 4.0 license	Open source means that the source code is freely available to all under GPL-3.0	Built by you means that the tool was developed by you or your organization
Tool ID	Standard Semantics	Standard Syntax
Tool ID is a unique identifier for each tool in the registry	Standard Semantics is a set of standard terms used to describe tools	Standard Syntax is a set of standard rules for describing tools
Community-driven	Backed by EDAM	Tool Platform
Community-driven means that the tool is developed and maintained by a community of users	Backed by EDAM means that the tool is supported by the EDAM project	Tool Platform is a set of services that support the tool
API	Documentation	Support
API is a set of protocols for interacting with the tool	Documentation is a set of instructions for using the tool	Support is a set of services that help users with the tool

bio.tools principles

- Open data
 - Content is freely available to all under CC BY 4.0 license
- Open source
 - Source code is freely available to all under GPL-3.0
- Built by the community
 - 4800+ (and growing!) contributors
- Persistent IDs
 - Unique, persistent, human-readable resource identifiers
- Standard semantics
 - Scientific function of bio.tools resources can be precisely annotated in defined terms from the EDAM ontology, including common topics, operations, types of data and data formats
- Standard syntax
 - Resources adhere to a rigorous syntax
 - ~50 key scientific, technical and administrative attributes (4 required)
- Community-driven
- Backed by ELIXIR
 - bio.tools will remain free, open and maintained in the long-term

Data model behind: biotoolsSchema

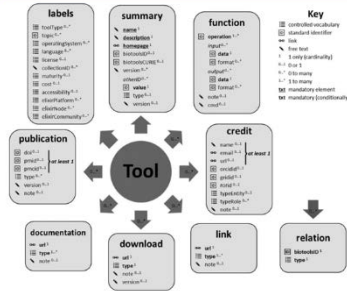
A simplified model which defined the attributes, information and scope provided to best describe software tools. A formal machine readable (and human understandable) schema to allow for data interoperability

- XML (JSON coming soon) schema
 - ~50 key scientific, technical and administrative attributes, uniform & rigorous syntax and semantics
- controlled vocabularies (18 in total)
 - e.g. tool type, software license, software maturity
- community-defined standard (v2.0, mature)
 - from multiple workshops/iterations
- compatible with related initiatives
 - shema.org/Bioschemas
- <https://doi.org/10.1093/gigascience/giaa157>

Data model: examples

- Required *
 - Name
 - bio.tools ID
 - Description
 - Homepage
- Labels:
 - Topics
 - Tool Type
 - OS
 - License
 - Language
 - Cost
- Publication
- Download links
- Documentation links
- Other links
- Credits
- Function

bio.tools provides [curation guidelines](#) for each attribute to help curators and regular users describe tools in a standard manner

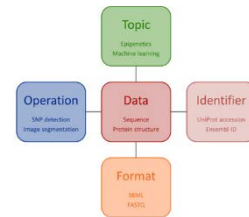


Scientific tool descriptions: EDAM ontology

Scientific function of bio.tools resources can be precisely annotated by concepts from the EDAM ontology

EDAM has 4 sections:

- Topic
- Operation
- Data (incl. Identifier)
- Format



EDAM ontology in bio.tools

- Topics in bio.tools define the scientific or technical domains the tool is developed for e.g. Transcription factors and regulatory sites, Genomics, Gene regulation
- Operations, Data and Formats are used in the context of functions
- Functions in bio.tools are represented as Input / Operation / Output triplets



EDAM has a hierarchical* structure

Topics

Operations

Browse it via <https://edamontology.github.io/edam-browser>

- Biology
 - Biomedical science
 - Chemistry
 - Computational biology
 - Biomolecular simulation
 - Function analysis
 - Molecular genetics
 - Chromogenetics
 - Gene and protein families
 - Gene expression
 - Gene structure
 - Functional, regulatory and non-coding RNA
 - Gene transcripts
 - Mobile genetic elements
 - Genetic variation
 - DNA mutation
 - DNA polymorphism
 - Structural variation
 - Molecular interactions, pathways and networks
 - Informatics
 - Phylogeny
 - Proteins
 - Sequence analysis
 - Sequence sites, features and motifs
 - Structure analysis
 - Computer science
 - Experimental design and studies
 - Informatics
 - Laboratory techniques
 - Laboratory and language
 - Mathematics
 - Medicine
 - Physics
 - Alignment
 - Fold recognition
 - Sequence alignment
 - Global alignment
 - Local alignment
 - Multiple sequence alignment
 - Pairwise sequence alignment
 - Sequence profile alignment
 - Structure-based sequence alignment
 - Tree-based sequence alignment
 - Structure alignment
 - Analysis
 - Annotation
 - Image annotation
 - Phylogenetic tree annotation
 - Sequence annotation
 - Sequence tag mapping
 - Text annotation
 - Calculation
 - Classification
 - Clustering
 - Comparison
 - Conversion
 - Correlation
 - Data handling
 - Design
 - Dimensionality
 - Indexing
 - Mapping
 - Modelling and simulation
 - Optimization and refinement
 - Prediction and recognition
 - Quantification

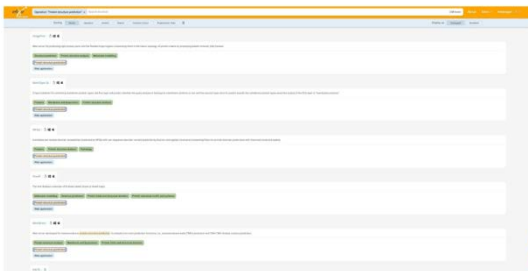
Collaboration with other projects

- EDAM Ontology: <http://edamontology.org>
- DebianMed: <https://www.debian.org/devel/debian-med>
- Galaxy: <https://usegalaxy.eu>
- BioContainers: <http://biocontainers.pro>
- Bioconda: <https://bioconda.github.io>
- SciCrunch: <https://scicrunch.org>
- ELIXIR TeSS: <https://tess.elixir-europe.org>
- FAIRSharing: <https://fairsharing.org>
- EuropePMC: <https://europepmc.org>
- ELIXIR Scientific Communities: <https://elixir-europe.org/communities>
- Australian BioCommons: <https://www.biocommons.org.au>
- MathWorks: <https://www.mathworks.com>
- OpenEBench: <https://openbench.bsc.es>
- Others

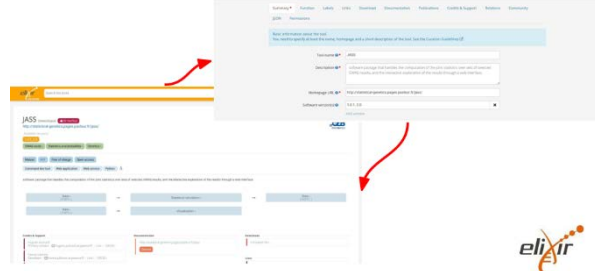
bio.tools useful links

- APIs
 - https://biotools.readthedocs.io/en/latest/api_reference.html
 - https://biotools.readthedocs.io/en/latest/api_usage_guide.html
- Documentation
 - <https://biotools.readthedocs.io>
 - <https://biotoolschema.readthedocs.io>
- Curation guidelines
 - https://biotools.readthedocs.io/en/latest/curation_guide.html
- Contact:
 - registry-support@elixir-mail.cbs.dtu.dk
 - <https://github.com/bio-tools/biotoolsRegistry/issues>
- EDAM:
 - <https://edamontology.org>
 - <https://github.com/edamontology/edamontology>

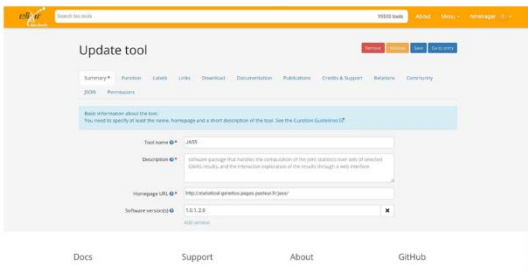
bio.tools: find tools



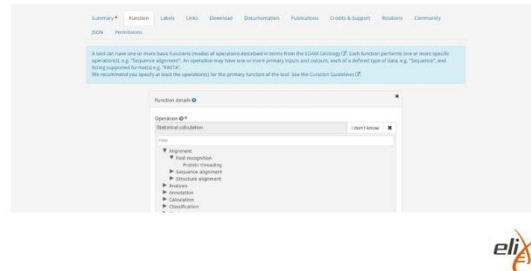
bio.tools: access and modify tool descriptions



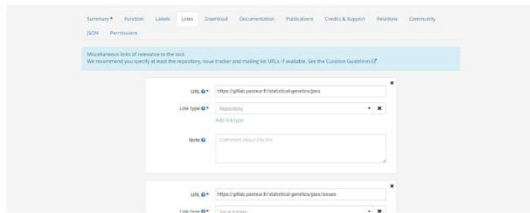
bio.tools curation: Summary



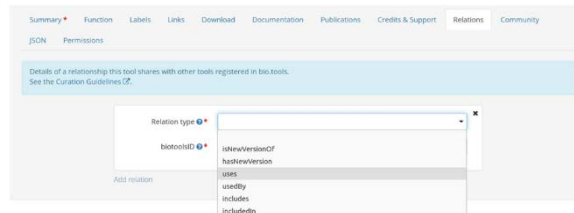
bio.tools curation: Function



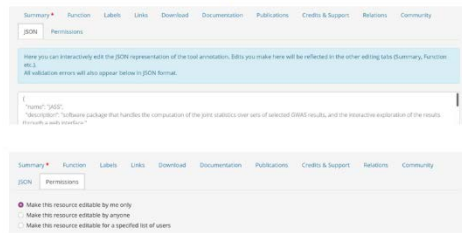
bio.tools curation: Labels, Links, etc.



bio.tools curation: Relations



bio.tools curation: JSON, Permissions



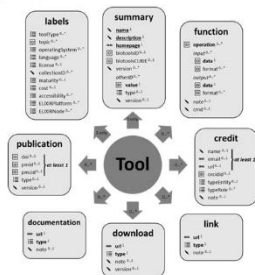
bio.tools descriptions

Add a new tool: required fields

Basic information about the tool. You need to specify at least the name, homepage and a short description of the tool. See the Curation Guidelines.

- Tool name *
- Tool description *
- Tool homepage URL *
- Tool identifier (biotoolsID) *:
 - Suggested from name (URL safe)
 - Editable if user wants
- Tool version (here or somewhere else)

biotoolsSchema



Tool properties links to docs

- [Accessibility](#)
- [Collection](#)
- [Cost](#)
- [Credit & Support](#)
- [Documentation](#)
- [Download](#)
- [ELIXIR community](#)
- [ELIXIR node](#)
- [ELIXIR platform](#)
- [Input data type](#)
- [Input file format](#)
- [License](#)
- [Links](#)
- [Maturity](#)
- [Operating system](#)
- [Operation](#)
- [Other ID](#)
- [Output data type](#)
- [Output file format](#)
- [Programming language](#)
- [Publication](#)
- [Tag as COVID-19](#)
- [Tool type](#)
- [Topic](#)

Labels

Miscellaneous scientific, technical and administrative details of the tool, expressed in terms from controlled vocabularies. We recommend you specify at least the tool type, license, and one or more topics. See the Curation Guidelines.

- Tag as COVID-19
- Tool type (Command-line, Web application, Desktop Application, Library)
- Operating system (Linux, Mac, Windows)
- Programming language (C++, Java, Python, R ...)
- Maturity (Emerging, Mature, Legacy)
- License (MIT, GPL-3.0, ..., Proprietary, Not licensed)
- Cost (Free, Free with restrictions, Commercial)
- Collection
- Accessibility (Open access, Open access (with restrictions), Restricted access)
- ELIXIR platform (Tools, Data, Interoperability, Compute, Training)
- ELIXIR node (Denmark, France, Germany, ...)
- ELIXIR community (Proteomics, Metabolomics, Rare diseases, ...)
- Tool confidence score (Tool, High, Medium, Low, Very Low)
- Other ID

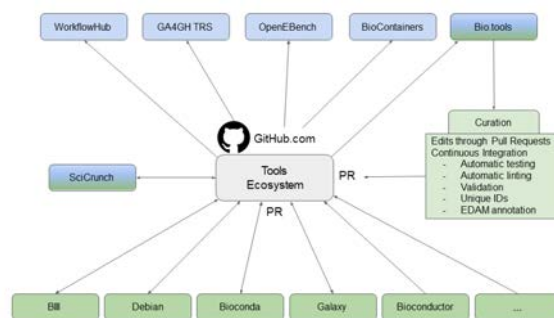
Links		Downloads	
Links <ul style="list-style-type: none"> • URL * • Linktype • Note 	Link types <ul style="list-style-type: none"> • Discussion forum • Galaxy service • Helpdesk • Issue tracker • Mailing list • Mirror • Repository • Service • Social media • Software catalogue • Technical monitoring • Other 	Downloads <ul style="list-style-type: none"> • URL * • Download type • Version • Note 	Download types <ul style="list-style-type: none"> • API specification • Binaries • Biological data • Command-line specification • Container file • Downloads page • Icon • Screenshot • Software package • Source code • Test data • Test script • Tool wrapper (CWL) • Tool wrapper (Galaxy) • Tool wrapper (Other) • Tool wrapper (Taverna) • VM image • Other

Documentation		Publication	
Documentation <ul style="list-style-type: none"> • URL * • Documentation type • Note 	Documentation types <ul style="list-style-type: none"> • API documentation • Citation instructions • Code of conduct • Command-line options • Contributions policy • FAQ • General • Governance • Installation instructions • Quick start guide • Release notes • Terms of use • Training material • User manual • Other 	<ul style="list-style-type: none"> • Digital Object ID (DOI) • PubMed ID • PubMed Central ID • Publication type <ul style="list-style-type: none"> ◦ Primary ◦ Method ◦ Usage ◦ Benchmarking study ◦ Review ◦ Other • Version • Note 	

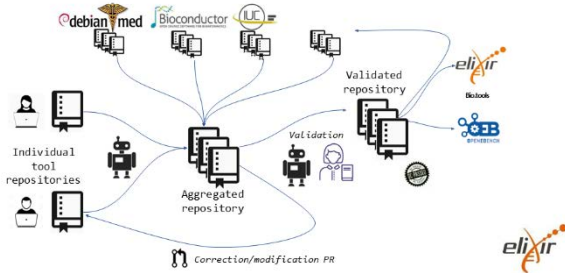
Credits & Support		Tool properties links to docs	
<ul style="list-style-type: none"> • ORCID ID • Name • Email • URL • gridid • rorid • fundrefid • Entity type • Entity role • Note 	Entity types <ul style="list-style-type: none"> • Person • Project • Division • Institute • Consortium • Funding agency Entity roles <ul style="list-style-type: none"> • Developer • Maintainer • Provider • Documentor • Contributor • Support • Primary contact 	<ul style="list-style-type: none"> • Accessibility • Collection • Cost • Credit & Support • Documentation • Download • ELIXIR community • ELIXIR node • ELIXIR platform • Input data type • Input file format • License 	<ul style="list-style-type: none"> • Links • Maturity • Operating system • Operation • Other ID • Output data type • Output file format • Programming language • Publication • Tag as COVID-19 • Tool type • Topic

“Tools Ecosystem”

(a work in progress)



The new Tools Ecosystem architecture (WIP)



EDAM tools

EDAM Tool Annotator and ontology browsers for EDAM

- https://bio.tools/static/eta_ (EDAM Tool Annotator)
- 3rd-party tooling to allow easier annotation of bio.tools entries with EDAM
- Provides an improved EDAM concept searching, and browsing over the bio.tools
- Steps:
 - Import tool metadata from bio.tools directly (via ID) or other sources
 - Annotate / edit EDAM-related tool metadata (topics, functions, ...)
 - Result is a bio.tools-compatible JSON
 - Request EDAM concept/term if none match your preferences
- Ontology browsers for EDAM:
 - https://bioportal.bioontology.org/ontologies/EDAM_ (NCBO BioPortal)
 - https://edamontology.github.io/edam-browser_ (EDAM Browser)
 - https://www.ebi.ac.uk/ols/ontologies/edam_ (OLS ontology browser)

EDAM concepts from text mining

- EDAMmap
 - Obtain EDAM concepts from text mining (free text, publications, etc.)
 - <https://bit.cs.ut.ee/edammap/>
 - <https://github.com/edamontology/edammap>
- Pub2tools
 - Generate bio.tools-compatible JSON data (including EDAM) from text mining publications
 - <https://github.com/bio-tools/pub2tools>

Demo 

Anexo C. Seminario de encapsulado de componentes software en contenedores

IMPACT-Data. Workshop en Contenerización de Software

<https://bit.ly/36HH1oo>

elixir

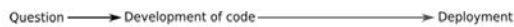
Björn Grüning
Bioconda/BioContainers community

www.elixir-europe.org

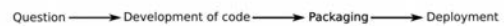
Disclaimer

This talk is titled with Containers and in the end you will get containers, but most of the talk will be about conda - optimized Containers and much more you get for free :)

Tool deployment & sustainability in science



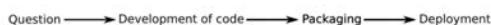
Tool deployment & sustainability in science



Package managers are charged with the task of finding, installing, maintaining or uninstalling software packages upon the user's command.

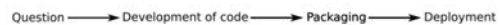


Tool deployment & sustainability in science



Tool deployment & sustainability in science

No standard available



What is needed?

- Programming language agnostic
- OS independent
- No root privileges needed
- Management of multiple version
- HPC and Cloud compatible
- easy to maintain



CONDA

- Open Source package manager
- Independent of any programming language and OS
- Fast, robust and easy package installation
 - > `conda install deeptools`
- Installation and management of multiple versions

CONDA

```

package:
  name: deeptools
  version: 3.11.0

source:
  file: deeptools-3.11.0.tar.gz
  url: https://pypi.python.org/packages/c/4/4/7/443b0d02205d0f11c3e077017f420c0522676760dee76b1e-3.11.0.tar.gz
  sha256: 17104b0d02205d0f11c3e077017f420c0522676760dee76b1e-3.11.0

requirements:
  build:
    - python
    - setuptools
    - numpy >=1.8.0
    - scipy >=0.17.0
    - matplotlib >=1.4.0
    - pyyaml >=4.2
    - pyzmq >=9.2.0
    - qt
  run:
    - python
    - pyzmq >=9.2.0
    - numpy >=1.8.0
    - scipy >=0.17.0
    - matplotlib >=1.4.0
    - pyyaml >=4.2
    - pyzmq >=9.2.0
  test:
    imports:
      - deeptools
  console_scripts:
    - bamCompare --version

about:
  home: https://github.com/eli4ir/deeptools
  license: GPL
  summary: A set of user-friendly tools for normalization and visualization of deep-sequencing data
    
```



CONDA

- Build packages
 - > `conda build packages/deeptools`
- building a scientific community

BIOCONDA
- using a unified build environment
 - > `circelci build`
- joining other Conda communities

<https://bioconda.github.io> | <https://conda.io/docs>



BIOCONDA

Navigation
 User Docs
 Contributing to Bioconda
 Developer Docs
 Bioconda on GitHub
 Package Index
 Quick search

BIOCONDA

Bioconda is a channel for the conda package manager specializing in bioinformatics software. Bioconda consists of:

- a repository of recipes hosted on GitHub
- a build system turning these recipes into conda packages
- a repository of packages containing over 7000 bioinformatics packages ready to use with `conda install`
- over 800 contributors and 570 members who add, modify, update and maintain the recipes

The conda package manager makes installing software a vastly more streamlined process. Conda is a combination of other package managers you may have encountered, such as pip, CPAN, CRAN, Bioconductor, apt-get, and homebrew. Conda is both language- and OS-agnostic, and can be used to install C/C++, Fortran, Go, R, Python, Java etc programs on Linux, Mac OSx, and Windows.

Conda allows separation of packages into repositories, or channels. The main default channel has a large number of common packages. Users can add additional channels from which to install software packages not available in the default channel. Bioconda is one such channel specializing in bioinformatics software.

Browse packages in the Bioconda channel: [Package Index](#)

Each package added to Bioconda also has a corresponding Docker BioContainer automatically created and uploaded to Quay.io. A list of these and other containers can be found at the [Biocontainers Registry](#)

<https://bioconda.github.io>
<https://conda.io/docs>

BIOCONDA

BIOCONDA

Navigation
 Contributing to Bioconda
 Developer Docs
 Package Index
 Quick search

Contributing to Bioconda

Bioconda is completely dependent on contributors to write, update, and maintain recipes. Every user will find below are instructions for one-time setup as well as a general procedure to follow for each recipe you like to add.

The basic workflow is:

- Follow the [recipe setup](#) instructions to get a local copy of the recipe repository
- Write a recipe or modify an existing one. A recipe consists of a metafile (the one appearing in this page) to install & test an example Python package
- Push
- Push your changes to GitHub. This triggers automatic building and testing of the recipe.
- Once the tests pass, the recipe is accepted by other members and their merged into the master branch. The resulting conda package and channels that are built on the master branch are uploaded to public repositories for searching use.

There are some details to be aware of, and some software can be challenging to package. The topics below provide more details.

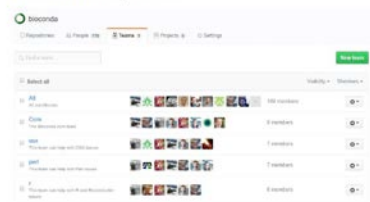
- recipe setup
- Contributor workflow
- Testing Recipes Locally
- Troubleshooting
- Build system
- Guidelines for Bioconda recipes
- LICENSE
- Pypi
- Uploading recipes
- Uploading recipes for a pending change
- Conda build 2.1

<https://bioconda.github.io>
<https://conda.io/docs>



BIOCONDA

- 9218 package (2022.04)
- contribute at <https://github.com/bioconda/bioconda-recipes>
- community driven



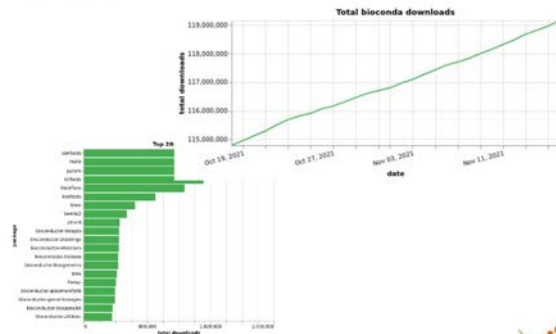
Community



1.419 Contributors ~30.000 merged PR



A few stats ...



Bioconda website:

<https://bioconda.github.io/environments/deepools/2021-04-16.html>



But the new cool kid is called Containers



But the new cool kid is called Containers

Question → Development of code → Packaging → Deployment



But the new cool kid is called Containers

Question → Development of code → Packaging → Deployment

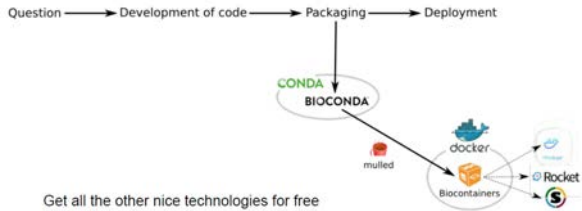
Mulled

- A layer donning approach to build containers
- without Dockerfile
- completely automatized
- powered by CI

Sharing the same artefact automatically, across technologies.



Container build without Dockerfile



- Get all the other nice technologies for free
- rkt
 - Singularity
 - Docker
 - ~40.000 containers for ~8.000 tools



Automated Container builds

package:
samtools: 1.3.1

vs.

```
process amberLeap {
  container="quay.io/biocontainers/samtools:1.3.1"
  input:
}
```

> mulled-build build-and-test 'samtools=1.3.1' --test 'samtools --help'

docker pull quay.io/biocontainers/samtools:1.3.1



Automated Container builds

package:
samtools: 1.3.1
bedtools: 2.26

> mulled-build build 'samtools=1.3.1,bedtools=2.26'

Automated Container builds

package:
samtools: 1.3.1
bedtools: 2.26

> mulled-build build 'samtools=1.3.1,bedtools=2.26'



Are multi-tool containers really needed? Decompose your workflows!

How do you name them?



No pets anymore, just cattles.



Findability!

Findability

- no search needed, just retrieve or fail
- build containers on the fly from packages
- build containers in-advance by monitoring GitHub repos

Automated Container builds

package:

samtools: 1.3.1

bedtools: 2.26

- predictable namespace for <samtools + bedtools> ?
- normalize package names, hash them
- normalize versions, hash them

```
mulled-v2-619c3451acc46ef686f3602375e74fb16bc237935232e184b5187
```

- <https://github.com/BioContainers/multi-package-containers>



Create containers via single line PRs

#Targets	base_image	image_build
#bedtools_1.1_Samtools-1.3.1+bedtools-2.26	quay.io/centos/centos7	0



Let a bot do it for you!

Add container mulled-v2-8186960447c5cb2faa697666dc1e6d919ad23f3e:5dde166a2021da0fda1facatc818a1ea0ff53ac7. #1286

```
Dockerfile
FROM quay.io/centos/centos7
RUN yum install -y samtools bedtools
```

Bots

- new software is released
- bot → conda
- community review
- bot → container
- (bot → missing container)
- (bot → missing singularity container)
- bot → Galaxy tool
- community review
- (bot → multi-package container)
- bot → workflow testing



Who is using it?



links

- cheat-sheet:
<https://raw.githubusercontent.com/bioconda/bioconda-outreach/master/SMB2019/CheatSheet.pdf>
- binder online bash:
<https://mybinder.org/v2/gh/qjbex/BinderBash/master>
- https://drive.google.com/file/d/1R256MvMY_j_OGKzN7e8daVMPUbuEsiX6/view?usp=sharing

Acknowledgements



Anexo D. Seminario Workflows - Uso de gestores de workflows, registro y compartición

Nextflow and nf-core IMPACT-Data workshop 30th May 2022

Cedric Notredame's lab
Jose Espinosa-Carrasco



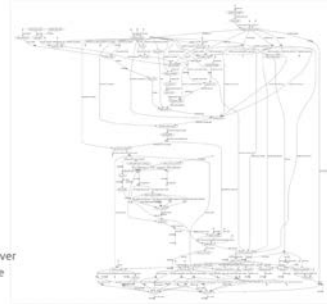
Bioinformatics workflows

- Data analysis applications to extract information from (large) genomic datasets.
- Embarrassingly parallelisation, can spawn 100s-100k jobs over a distributed cluster.
- Mash-up of many different tools and scripts (dependencies).
- Complex dependency trees and configuration

Very fragile ecosystem

The same pipeline deployed in different environments produces different results (!)

An example



Steinbiss et al., Companion: a web server for annotation and analysis of parasite genomes, DOI: 10.1093/nar/gkw292

Workflow Management System to the rescue

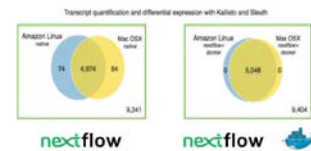
Published: 11 April 2022
Nextflow enables reproducible computational workflows

Paolo Di Tommaso, Maria Chabalina, Evan H. Fodor, Pablo F. del Barrio, Emilio Palumbo & Cedric Notredame

Nature Methods 19, 316–319 (2022) | [View this article](#)



Paolo Di Tommaso, Nextflow Lead



Workflow managers enable reproducibility

Reproducible, scalable, and shareable analysis pipelines with bioinformatics workflow managers

Laura Wratzen, Andreas Wilm & Jonathan Göke

Nature Methods 18, 1161–1168 (2021) | Cite this article

The rapid growth of high-throughput technologies has transformed biomedical research. With the increasing amount and complexity of data, scalability and reproducibility have become essential not just for experiments, but also for computational analysis. However, transforming data into information involves running a large number of tools, optimizing parameters, and integrating dynamically changing reference data. Workflow managers were developed in response to such challenges. They simplify pipeline development, optimize resource usage, handle software installation and versions, and run on different compute platforms, enabling workflow portability and sharing. In this Perspective, we highlight key features of workflow managers, compare commonly used approaches for bioinformatics workflows, and provide a guide for computational and noncomputational users. We outline community-curated pipeline initiatives that enable novice and experienced users to perform complex, best-practice analyses without having to manually assemble workflows. In sum, we illustrate how workflow managers contribute to making computational analysis in biomedical research shareable, scalable, and reproducible.

nextflow main features



Task example

```
process align_sample {
    input:
    path ref_fasta
    path sample_fastq

    output:
    path 'sample.bam', emit: bam

    script:
    """
    bwa mem $ref_fasta $sample_fastq \
    | samtools sort -o sample.bam
    """
}
```

Dataflow paradigm

- Declarative computational model for parallel process executions
- Processes wait for data, when an input set is ready the process is executed
- They communicate by using dataflow variables i.e. async FIFO queues called channels
- Parallelisation and tasks dependencies are implicitly defined by process in/out declarations



Comparison with other workflow managers

Perspective | Published: 23 September 2021

Reproducible, scalable, and shareable analysis pipelines with bioinformatics workflow managers

Laura Wratzen, Andreas Wilm & Jonathan Göke

Nature Methods 18, 1161–1168 (2021) | Cite this article

Table 1 | Overview of workflow managers for bioinformatics. This article compares features across several workflow managers.

Tool	Flow	Based on	Experiment	Portability	Scalability	Sharing	Pipeline
Galaxy	Graphical	***	***	***	***	***	***
Nextflow	Scripting	***	***	***	***	***	***
Cromwell	DSL	***	***	***	***	***	***
Snakemake	DSL	***	***	***	***	***	***
WDL	DSL	***	***	***	***	***	***
Apache Airflow	DSL	***	***	***	***	***	***
Luigi	Library	***	***	***	***	***	***
OpenWDL	Scripting	***	***	***	***	***	***
WDL + Cromwell	Scripting	***	***	***	***	***	***
WDL + Cromwell	Scripting	***	***	***	***	***	***

Task example

```
bwa mem reference.fa sample.fq \
| samtools sort -o sample.bam
```

Task composition

```
process align_sample {
    input:
    path ref_fasta
    path sample_fastq

    output:
    path 'sample.bam', emit: bam

    script:
    """
    bwa mem reference.fa $sample_fastq \
    | samtools sort -o sample.bam
    """
}

process index_sample {
    input:
    path bam_to_idx

    output:
    path "${bam_to_idx}.bai", emit: bai

    script:
    """
    samtools index $bam_to_idx
    """
}

workflow {
    align_sample ( ref_fasta_ch, fastq_ch )
    index_sample ( align_sample.out.bam )
}
```

How parallelization works

```
samples_ch = Channel.fromPath("data/sample.fastq")

process FASTQC {
    input:
    path reads
    output:
    path 'fastqc_logs', emit: fastqc_ch

    """
    mkdir fastqc_logs
    fastqc -q ${reads} -f fastqc -o fastqc_logs
    """
}

workflow {
    FASTQC ( samples_ch )
}
```

How parallelization works

```

samples_ch = Channel.fromPath("data/*fastq")

process FASTQC {
  input:
  path reads
  output:
  path 'fastqc_logs', emit: fastqc_ch
  ...
  mkdir fastqc_logs
  fastqc -q ${reads} -f fastq -o fastqc_logs
  ...
}

workflow {
  FASTQC ( samples_ch )
}
    
```

Supported platforms



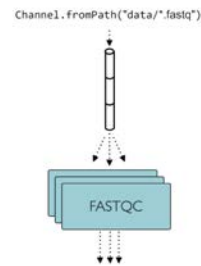
Portability



Portability



Implicit parallelism



Portability



```
nextflow run your-script.nf --with-geo@geopkg/nextflow
```

Portability



Portability



Containerisation



• Nextflow envisioned the use of software containers to fix computational reproducibility

• Mar 2014 (ver 0.7), support for Docker

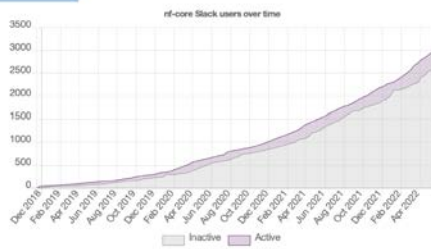
• Dec 2016 (ver 0.23), support for Singularity



What makes nextflow strong???



Slack members over time



Response time to pull requests and issues



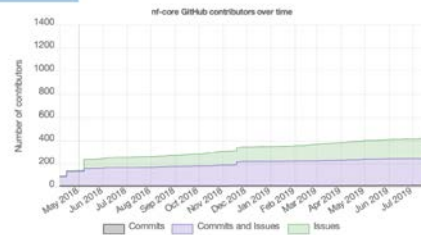
Nextflow turns modular

- DSL-2 is an extension of Nextflow syntax (released on version 20.07.1)
- Enables the definition of reusable modules and sub-workflows
- Pipeline can be now split in different files instead of having a huge script with all the logic
- Some changes to respect DSL-1 but don't change the core Nextflow concepts (e.g. channels can be reused without need of create multiple copies of the same channel)
- DSL2 became the default mode in the last stable release 22.04.1

GitHub nf-core members over time



nf-core GitHub contributors



nf-core community resources

<https://nf-core.re/>

<https://github.com/nf-core/>

<https://nf-co.re/join/slack>

https://twitter.com/nf_core

<https://groups.google.com/forum/#!forum/nf-core>

<https://www.youtube.com/c/nf-core>

nf-core/bytesize

All Tuesdays 13:00 CEST

Bytesize: resources to learn Nextflow

What nf-core is?

- A community of users and developers
- A curated set of analysis pipelines build using Nextflow
- A set of guidelines (standard)
- Helper tools



nf-core guidelines

Build using Nextflow
 Have a MIT license
 Software bundled using Docker/Singularity
 CI testing (include minimal test dataset)
 Pass of core list tests
 Stable release tags (GitHub releases with associated changelog, and DOI using Zenodo)
 Common pipeline structure and usage
 Run in a single command (pipelines should not be split into multiple sub-pipelines)
 Comprehensive documentation
 A responsible point of contact

Software bundled with Bioconda (package the software using Bioconda if not already available)
 Explicit support for cloud environment
 Benchmarks from running on cloud environments
 Optimized output file formats (Standard format where possible, including CRAM)

nf-core requirements

nf-core requirements

The nf-core framework for community-curated bioinformatics pipelines

nf-core recommendations

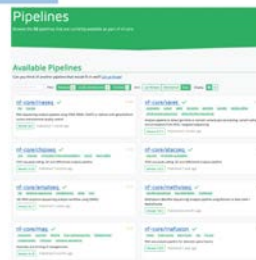
Why having strict guidelines

Pipelines work in a comparable manner

```
nextflow run <pipeline> -profile test,singularity
nextflow run <pipeline> -profile test,docker
nextflow run <pipeline> -profile test,conda
```

A common interface: e.g. web form

nf-core curated pipelines



- ✓ 36 RELEASED
- 🔧 24 UNDER DEVELOPMENT
- 📁 6 ARCHIVED

Why having strict guidelines

- Follow FAIR principles
- Adhere to current best practices in terms of computational reproducibility and interoperability
- Guarantee the portability between different computational infrastructures
- Enable a set of common features between pipelines (how they run, documentation, etc.)

A common interface: e.g. web form

A package of helper tools

nf-core/
tools

A python package with helper tools for the nf-core community



`pip install nf-core`

BIOCONDA

`conda install -c bioconda nf-core`



`docker pull nfcore/tools`

nf-core tools



nf-core DSL2 concepts

MODULE: A process that can be used within different pipelines and is as atomic as possible i.e. cannot be split into another module.
e.g. a module file containing the process definition for a single tool such as FastQC

SUB-WORKFLOW: A chain of multiple modules that offer a higher-level functionality within the context of a pipeline.
e.g. a sub-workflow to sort, index and run some basic stats on a BAM file.

WORKFLOW: An end-to-end pipeline created by a combination of Nextflow DSL2 individual modules and sub-workflows.
e.g. from one or more inputs to a series of final inputs

Nextflow tower

- Web user interface to interact with Nextflow
- An API to “talk” to pipelines
- Seamless configuration of cloud environments
- Enables to run pipelines in the cloud or HPC



Join the community

nextflow

<https://nextflow.io>

nf-core

<https://nf-co.re/>

<https://nf-core/join>

https://join.slack.com/t/nextflow/shared_invite/zt-11wltw5-R6SNBpVksOJAx5SPOXNrZg

nf-core turns DSL2

nf-core/maseq
nf-core/viralrecon
nf-core/fetchngs
nf-core/amplicon
nf-core/bacass
nf-core/fetchngs
nf-core/boellmagic
nf-core/cutandrun
nf-core/amplicon
nf-core/maseq
nf-core/mfusion
...



nf-core DSL2 modules

Modules

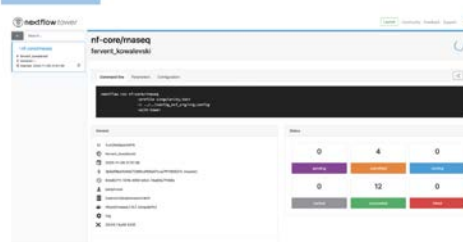
Browse the 311 modules that are currently available as part of nf-core.

Available Modules

Modules are the building blocks of all DSL2 nf-core blocks. You can find more info, if you would like to write your own module.

fastq	abacast	genome	genome	genome
fastq	genome	genome	genome	genome
reference	reference	reference	reference	reference
index	index	index	index	index
genome	genome	genome	genome	genome
fastq	fastq	fastq	fastq	fastq
fastq	fastq	fastq	fastq	fastq
assembly	assembly	assembly	assembly	assembly

Nextflow tower enables to launch and monitor your executions



Demo

Thanks!!!

nf-core <https://nf-core.net/>

nextflow <https://www.nextflow.io>

<https://www.nextflow.io/doc/latest/index.html>

nextflow tower <https://tower.nf>

Notredame Lab
Cedric Notredame

Pablo d Tommaso
Sequera CTO & co-founder

Evan Fioden
Sequera CEO & co-founder

Phil Ewels
SciLifeLab Sweden
nf-core team

PRIORITY OF CONFIGS

- Parameters specified on the command line (`--something value`)
- Parameters provided using the `--params-file` option
- Config file specified using the `-c my_config` option
- The config file named `nextflow.config` in the current directory
- The config file named `nextflow.config` in the workflow project directory
- The config file `$HOME/.nextflow/config`
- Values defined within the pipeline script itself (e.g. `main.nf`)

Example: use software bundle with Bioconda/Biocontainers

bioconda / packages / fastqc 0.11.9

A quality control tool for high throughput sequence data.

License: GPL v=3
Home: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
446619 total downloads
Last upload: 1 year and 15 days ago

BioContainers **docker**

```
conda (params.enable_conda ? "bioconda::fastqc@0.11.9" : null)
container "${ workflow.containerEngine == "singularity" && !task.ext.singularity_pull_docker_container ?
"https://depot.galaxyproject.org/singularity/fastqc@0.11.9-0" :
"quay.io/biocontainers/fastqc@0.11.9-0" }
```

HOW TO CONTRIBUTE

1. Join the nf-core community

#new-pipelines

There should only be a single pipeline for a given data + analysis type

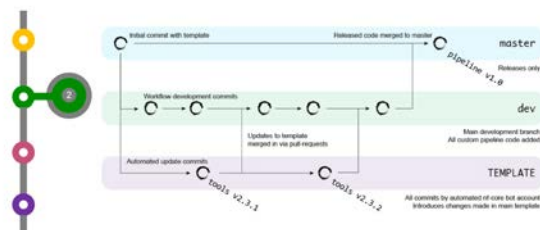
CREATE A PIPELINE FROM THE TEMPLATE

nf-core create

```
nf-core/tools, version 2.3.2 - https://nf-co.re
Description: de novo assembly pipeline
Author: Jose Espinosa-Carrasco
Usage: Creating new nf-core pipeline: "nf-core/assembly"
Initiating pipeline git repository
Done. Remember to add a remote and push to Github
git remote add origin git@github.com:USERNAME/REPO_NAME.git
git push --all --origin
This will also push your newly created dev branch and the TEMPLATE branch for syncing
!!!!!! SUCCESS!!!!!!
If you are interested in adding your pipeline to the nf-core community,
please open an issue on the nf-core Github repository and contact
the community.
```

Use the template!!

SOFTWARE DEVELOPMENT CYCLE



CONFIGURATION FILES

- Default config (e.g. base.config)
 - Automatically loaded
 - Sensible default resources request
- Core profiles (e.g. docker, singularity, test) → `-profile test,docker`
 - Specify software packaging
 - Specify common presets
- Institutional profiles (nf-core/configs) → `-profile CRG`
 - Specify job submission for your institution cluster
 - Specify software packaging and other settings
 - Available for all nf-core pipelines
- Your local config files → `-c flag e.g. -c my_local.config`
 - Custom resource requirements, user specific parameters

nextflow-schema.json

```
params.foo = "BovReg"
println ("Hello $params.foo !")

$ nextflow run main.nf --foo EuroFAANG Hello EuroFAANG !
```

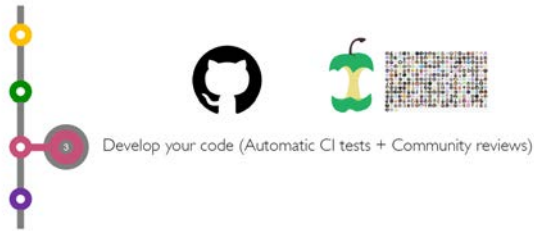
```
nextflow-schema.json
"foo": {
  "type": "string",
  "default": "BovReg",
  "description": "A greetings parameter"
}
```

Automates

- User input validation
- Pipeline CLI help
- Documentation
- User interfaces

nf-core schema build

HOW TO CONTRIBUTE



HOW TO CONTRIBUTE



7,8 cm

nf-core list

Pipeline Name	Stars	Latest Release	Released	Last Pulled	New latest release?
clonem	3	1.0.0	12 hours ago	-	-
gprn	2	1.0.0	15 hours ago	-	-
fluprotemics	5	1.2.3	yesterday	-	-
chipseq	72	1.2.2	5 days ago	-	-
snor	49	2.3.3	3 weeks ago	-	-
scRNAseq	19	1.1.0	4 weeks ago	-	-
methyseq	63	1.0	1 months ago	-	-
bedtools	3	1.0.0	2 months ago	-	-
mc	86	1.2.0	3 months ago	-	-
sm1199	59	1.2.0	3 months ago	-	-
snitch	59	2.2	3 months ago	-	-
cohort	3	1.0.2	3 months ago	-	-
chrseq	313	3.0	4 months ago	8 months ago	No (v1.4.2)

nf-core launch rnaseq

```

Select Nextflow flags to control how the pipeline runs.
These are not specific to the pipeline and will not be saved in any persistent file. They are just used when including the search comment.
(Use arrow keys)
--mode (v1.0)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

nf-core launch rnaseq

nf-core launch rnaseq

Define where the pipeline should find input data and save output data.

Input/output options

Input directory (required): /mnt/1000

Output directory (required): /mnt/1000

Reference genome options

Options for the reference genome indices used to align reads.

Reference genome (required): hg38

Reference genome index (required): hg38

Reference genome index (optional): hg38

Galaxy + WorkflowHub

Ignacio Eguinoa (VIB - ELIXIR BE)
30/05/2022



Agenda

- Galaxy
 - Overview
 - Tools
 - Templates
 - Ejemplos
 - Galaxy workflows
 - Editor
 - Ejecución
- WorkflowHub
 - Registrar un workflow
 - RO-Crate

Galaxy Project

Proporciona una capa de abstracción para acceder y hacer uso de la infraestructura de cómputo. Una simple interfaz gráfica provee una capa de abstracción para:

- Ejecutar comandos en entorno Linux.
- Definir y ejecutar un workflow.
- Resolver dependencias de ejecución.
- Job scheduling.
- Data management.
 - Metadata + tipos de datos.
- Manejo de usuarios.



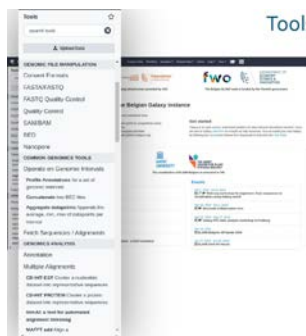
Galaxy es toda una plataforma/ecosistema de abstracciones y herramientas, pero el core se basa en 2 entidades: Herramientas y Workflows.



¿Cómo usar la plataforma?

- Instalación local:
 - Clone <https://github.com/galaxyproject/galaxy>
 - ./run.sh
 - Default configuration
- Versiones distribuidas en contenedores Docker:
 - <https://hub.docker.com/r/irbgruening/galaxy-stable/>
- Servidores públicos:
 - Usegalaxy.org
 - Usegalaxy.be
 - Usegalaxy.eu
 - Usegalaxy.es
 - Usegalaxy.org.au
 - ...

Tools: unidad atómica de ejecución.



- Listadas en el panel de herramientas (panel a la izquierda), organizadas en secciones.
- Módulos instalables: cada admin de un servidor decide qué herramientas instalar y poner a disponibilidad de usuarios.
- En Galaxy, "todo" es una herramienta: cargar datos locales, de repositorios remotos, enviar datos, procesamiento, linux commands, herramientas bioinformáticas, etc.
- Definidas mediante "tool wrappers" que proveen reproducibilidad.



GUI de cada herramienta



Tool wrappers: templates para cada herramienta.

Schema XML
<https://docs.galaxyproject.org/en/latest/dev/schema.html>

```
<tool id="hello" name="hello" version="0.01">
  <description>World's description</description>
  <command>-[CDATA[
    /bin/echo "Hello World"
  ]]></command>
  <inputs>
    <param name="mystring" type="text" label="Say something interesting"/>
  </inputs>
  <outputs>
    <data format="tabular" name="output1" label="hello_world"/>
  </outputs>
  <help>-[CDATA[
    **What it does**
  ]]></help>
</tool>
```



Reproducible job execution



¿Cómo definir las dependencias?

La sección de requirements del wrapper es un indicador a la plataforma sobre qué dependencias deben estar disponibles en el entorno de ejecución

<https://docs.galaxyproject.org/en/latest/dev/schema.html#tool-requirements>

La plataforma puede estar configurada para resolver dependencias usando distintos sistemas y prioridades (mayor o menor reproducibilidad):

- Paquetes locales (deprecated)
- Conda (el admin debe configurar los canales y prioridades)
- Contenedores (e.g. <https://quay.io/biocontainers/fastqc.0.11.2--1>)

Bioconda <-> Biocontainers

EDAM metadata

- <https://edamontology.org/>
 - <https://docs.galaxyproject.org/en/latest/dev/schema.html#tool-edam-topics>
- ```
<edam_operations>
 <edam_operation operation="3434" edam_operation="3434">
 </edam_operation>
 </edam_operations>
```
- Ayuda a enriquecer la metadata de la herramienta/proceso.
  - Se usan tanto dentro (e.g. para organizar el panel de herramientas) como en servicios externos (e.g. workflowhub)

## Repositorios + Cómo crear un wrapper para tu herramienta

Las herramientas son módulos instalables: wrapper + metadata

- Galaxy tiene su propio package manager interno que permite instalar paquetes desde <https://tools.bioinformatics.edu/>
- Cada uno puede subir sus propios paquetes (también en <https://galaxytoolshed.gli.gsc.psu.edu/>)
- Se recomienda usar Github/Gitlab para mantener el código del paquete.
- Galaxy tiene su propio repositorio (IUC => nf-core) el cual sigue las mejores prácticas de desarrollo + CI test + updates automáticos, etc

Muchas herramientas ya tienen un wrapper en el ecosistema Galaxy

Para hacer un wrapper desde cero:

- <https://training.galaxyproject.org/training-material/topics/dev/tutorial/tutorial-from-scratch/tutorial.html>
- Crear las dependencias en bioconda-biocontainers <https://github.com/bioconda/bioconda-recipes>
- El schema del XML (sistema de tipos de datos, outputs, etc): <https://docs.galaxyproject.org/en/latest/dev/schema.html>
- Usar los repositorios publicos de Galaxy Project (IUC) como guía: <https://github.com/galaxyproject/tools-iuc/tree/master/tools>

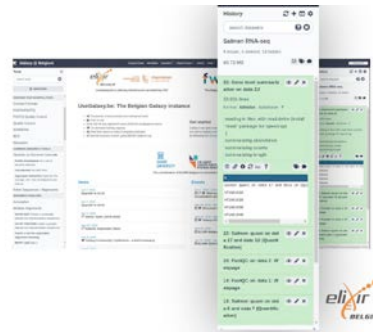


## Communities sharing one coherent framework

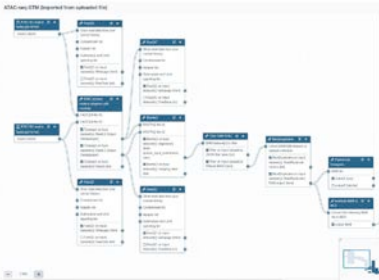
- [ma.usegalaxy.eu](https://ma.usegalaxy.eu)
- [clipseq.usegalaxy.eu](https://clipseq.usegalaxy.eu)
- [metagenomics.usegalaxy.eu](https://metagenomics.usegalaxy.eu)
- [hicexplorer.usegalaxy.eu](https://hicexplorer.usegalaxy.eu)
- [cheminformatics.usegalaxy.eu](https://cheminformatics.usegalaxy.eu)
- [proteomics.usegalaxy.eu](https://proteomics.usegalaxy.eu)
- [imaging.usegalaxy.eu](https://imaging.usegalaxy.eu)
- [metabolomics.usegalaxy.eu](https://metabolomics.usegalaxy.eu)
- [ecology.usegalaxy.eu](https://ecology.usegalaxy.eu)
- [nanopore.usegalaxy.eu](https://nanopore.usegalaxy.eu)
- [singlecellomics.usegalaxy.eu](https://singlecellomics.usegalaxy.eu)
- [humancellatlas.usegalaxy.eu](https://humancellatlas.usegalaxy.eu)
- [virology.usegalaxy.eu](https://virology.usegalaxy.eu)
- [climate.usegalaxy.eu](https://climate.usegalaxy.eu)
- [streetscience.usegalaxy.eu](https://streetscience.usegalaxy.eu)



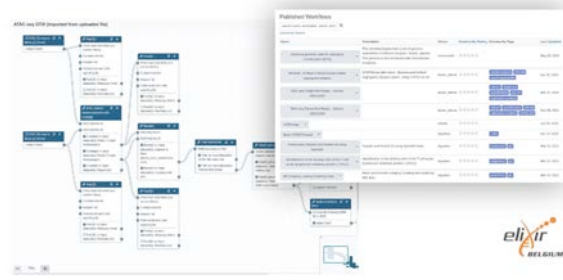
## Historial del análisis



## Graphical workflow editor



## Compartir workflows dentro de una misma instancia



## Formato de serialización de workflows

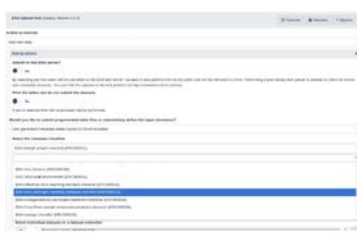
- Generado automáticamente a partir del editor de workflow.
- JSON based (formato v1)
  - YAML based (formato v2)
- Referencias al id de la herramienta -> no todos los detalles son exportados.
- Compartir inter-instancias de Galaxy: debe tener las mismas herramientas instaladas.

```

graph LR
 subgraph "Workflow 1"
 W1[Workflow 1]
 end
 subgraph "Workflow 2"
 W2[Workflow 2]
 end
 W1 --> W2

```

## Submission to public repositories



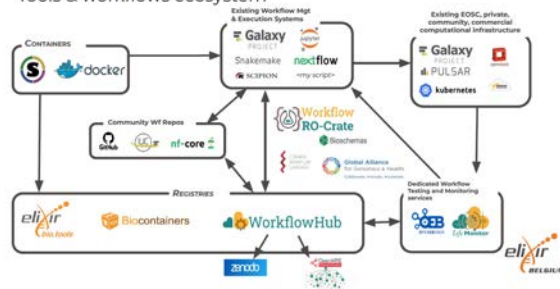
- Simplified submission interface.
- Multiple metadata input options.
- Include the submission automatically in your analysis.



## Demo

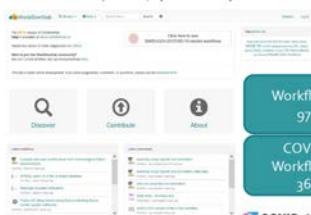
- Workflow input datasets
- Parámetros estáticos: <https://training.galaxyproject.org/training-material/topics/galaxy-interface/tutorials/workflow-parameters/tutorial.html>
- Descargar workflows
- Ejecutar workflows

## Tools & workflows ecosystem



## WorkflowHub.eu: a FAIR workflow registry

Beta release Sept 2020, sponsored by EOSC-Life



- Perpetual Development in the open by an open community
- Registering on behalf of makers
- Regular releases of new features
- Agile revisions of features

COVID-19 Data Portal



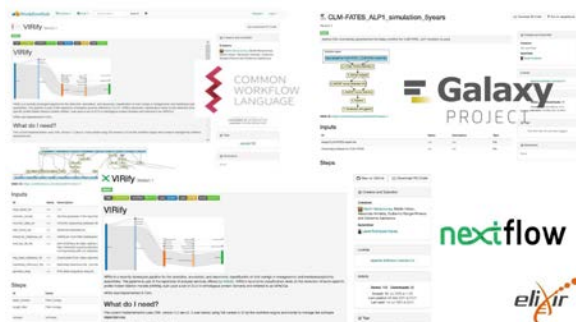
Workflows at the centre, as entry point for end-users

## Open to any workflow and WfMS

- Open to workflows from all disciplines and any country, currently already spanning
  - Life Science
  - Biodiversity
  - Climate change
  - Industrial Biotech
- Workflow management system agnostic
- Workflows may remain in their native repositories in their native form
- Registry & repository functionality







## Representación interna de un Workflow en Workflowhub

- Internamente WorkflowHub intenta llevar todos los workflows a un único estándar: CWL Abstract, con el objetivo de aumentar la interoperabilidad.
  - Lista los pasos y orden de ejecución
  - Describe inputs/outputs
  - No especifica cómo generar el comando que cada paso debe ejecutar (orden de parámetros, lógica, etc)
- Esta representación estándar permite(-iría):
  - Comparar la estructura de workflows independientemente de su lenguaje inicial.
  - Buscar un workflow usando una estructura abstracta o metadata (e.g procesos EDAM, inputs, tags, ...) y obtener la implementación en distintos lenguajes.

## Community driven standards to describe workflows and metadata

**COMMON WORKFLOW LANGUAGE**  
Canonical workflow description  
Native or Abstract CWL  
<https://www.common-workflow-language.org/>

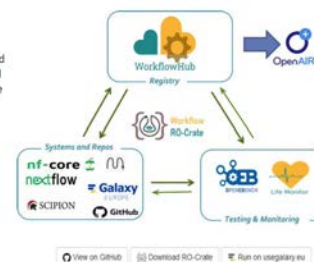
**Bioschemas**  
Metadata for registration and discovery  
ComputationalWorkflow and FormalParameter  
Schema.org profile and types  
<https://bioschemas.org/profiles/ComputationalWorkflow/>

**Global Alliance for Genomics & Health**  
API Exchanging Tools and Workflows  
GA4GH Tool Registry Service (TRS) API

**Workflow RO-Crate**  
Exchange format for WorkflowHub  
RO-Crate contains files or references to files in native repository, documentation, tests, examples

## Integration with Other Services and WfMS platforms

- Register a workflow directly as uploaded RO-Crate
- Register a workflow file, abstract CWL, and diagram, provide additional metadata and license information, and have an RO-Crate made for you
- Download an RO-Crate



### Coming soon

- API method for submitting RO-Crates to WorkflowHub.
- GitHub action to manage a repository contained workflows.

## Representación externa: RO-Crate (Research Object Crate)

A method of organizing file-based data with associated metadata, using linked data principles, in both human and machine readable formats, with the ability to include additional domain-specific metadata.

- The core of RO-Crate is a JSON-LD file, the RO-Crate Metadata File, named `ro-crate-metadata.json`. This file contains structured metadata about the dataset as a whole (the Root Data Entity) and, optionally, about some or all of its files.
- Schema.org is the base metadata standard for RO-Crate.
- Schema.org was chosen because it is widely used on the World Wide Web and supported by search engines, on the assumption that discovery is likely to be maximised if search engines index the content.
- <https://www.researchobject.org/ro-crate/1.1/context.jsonld>



## Workflow specific schema

### Schema.org es extensible:

- Thing > CreativeWork > SoftwareSourceCode > ComputationalWorkflow
- <https://bioschemas.org/profiles/ComputationalWorkflow/1.0-RELEASE>
  - Minimally needs to describe: author, input, output, version...



```

streetAddress
Thing > Property > streetAddress
The street address. For example, 1000 Amphitheatre Parkway.

Values expected to be one of these types
Text

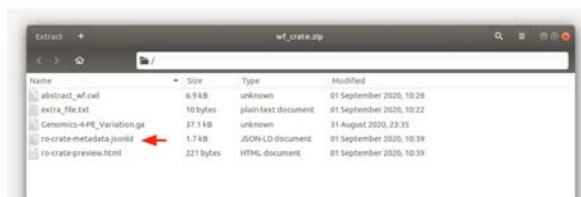
Used on these types
PostalAddress

Examples

Example 1
@context: "https://schema.org",
"@type": "Person",
"address": {
 "@type": "PostalAddress",
 "addressLocality": "Seattle",
 "addressRegion": "WA",
 "postalCode": "98002",
 "streetAddress": "20341 Whitworth Institute 405 N. Whitworth"
}

```

## Ejemplo de workflow RO-Crate (\*.zip)



```

{
 "name": "Genomics - PE Variant",
 "description": "Genomics - PE Variant",
 "author": "Genomics - PE Variant",
 "version": "1.0.0",
 "license": "MIT",
 "keywords": [
 "genomics",
 "variant calling",
 "bioinformatics"
],
 "dependencies": {
 "samtools": "1.10.0",
 "bcftools": "1.10.0",
 "pysnpSites": "1.0.0"
 },
 "entryPoint": "main.py",
 "mainFunction": "main",
 "mainModule": "main",
 "mainPackage": "main",
 "mainClass": "main",
 "mainMethod": "main",
 "mainFunctionName": "main",
 "mainModulePath": "main",
 "mainPackagePath": "main",
 "mainClassPath": "main",
 "mainMethodPath": "main"
}

```

General metadata

ro-crate-metadata.json

Workflow metadata



## What was done + work in progress

- Python implementation of RO-Crate: <https://github.com/ResearchObject/ro-crate-py>
- Entry point to extract workflow details from Galaxy: <https://github.com/galaxyproject/galaxy/pull/9407>
- Use the CWLProv profile (<https://github.com/common-workflow-language/cwlprov>) to export the provenance of a workflow run:
  - intermediate files + logs + generatedBy + executionTime + ...
  - Based on general provenance ontology (<https://www.w3.org/2013/05/trace-prov-20130430/>)
- Similar to: <https://academic.oup.com/gigascience/article/8/11/giz085/5611001>



## ¿Cómo registrar tu workflow?

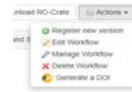
- Archivo principal del workflow:
- Upload de archivos individuales
  - Link a repositorios Git
  - Upload de RO-Crates
- Metadata:
- Completar formulario de upload
  - RO-Crate



## Handling Workflow Lifecycles

RO-Crates upload is a starting point

- Enhancing**
- Need to include references to **Datasets and Publications** associated with the workflow in WorkflowHub
  - Update and enhance uploaded RO-Crates with metadata provided by the WorkflowHub
  - RO-Crate download is enriched
- Versioning of workflow files or RO-Crate**
- Register a new version through the web interface
  - Control over visibility of past versions (unless they have a DOI, in which case they are public)



- Publishing**
- Generate DOI for a public workflow entry
  - Simple one-click to create
  - DOI is tied and resolves to a specific version
  - Citation is available to copy



## WorkflowHub - usegalaxy.eu integration

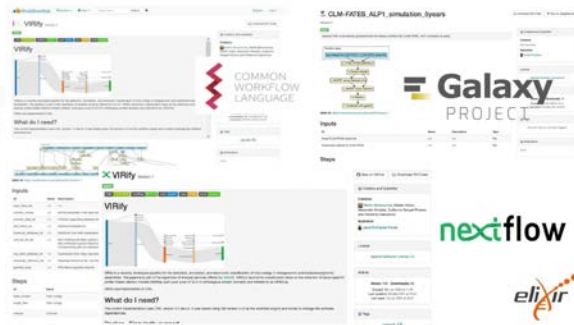
NOT an execution platform, but can be coupled to execution platform.



## Ejecutar workflows de WorkflowHub en Galaxy.



## Ejecutar workflows de WorkflowHub en Galaxy.



## Join the WorkflowHub Club

We gratefully acknowledge the WorkflowHub Club, Bioschemas Group, RO-Crate Group, CWL Community and our WFMS partners in Galaxy, Snakemake, Nextflow, CWL, SCIPION.

- **WorkflowHub Club**  
<https://about.workflowhub.eu>
- **Bioschemas**  
<https://bioschemas.org/groups/Workflow/>
- **Research Object/ RO-Crate**  
<https://www.researchobject.org/>
- **Common Workflow Language**  
<http://commonwl.org>

<https://workflowhub.eu>



## Acknowledgements

