

**BỘ GIÁO DỤC
VÀ ĐÀO TẠO**

**VIỆN HÀN LÂM
KHOA HỌC VÀ CÔNG NGHỆ VN**

HỌC VIỆN KHOA HỌC VÀ CÔNG NGHỆ



Đào Quang Hà

**ỨNG DỤNG TIN SINH HỌC TRONG PHÂN TÍCH
HỆ PHIÊN MÃ CÂY SÂM NGỌC LINH
(*Panax vietnamensis* Ha et Grushv.) 4 NĂM TUỔI**

**LUẬN VĂN THẠC SĨ
SINH HỌC THỰC NGHIỆM**

ĐÀO QUANG HÀ

**SINH HỌC
THỰC NGHIỆM**

2022

Hà Nội - 2022

**BỘ GIÁO DỤC
VÀ ĐÀO TẠO**

**VIỆN HÀN LÂM
KHOA HỌC VÀ CÔNG NGHỆ VN**

HỌC VIỆN KHOA HỌC VÀ CÔNG NGHỆ



Đào Quang Hà

**ỨNG DỤNG TIN SINH HỌC TRONG PHÂN TÍCH
HỆ PHIÊN MÃ CÂY SÂM NGỌC LINH
(*Panax vietnamensis* Ha et Grushv.) 4 NĂM TUỔI**

Chuyên ngành: Sinh học thực nghiệm
Mã số: 8420114

**LUẬN VĂN THẠC SĨ
NGÀNH SINH HỌC**

NGƯỜI HƯỚNG DẪN KHOA HỌC:

1. TS. Nguyễn Tường Vân
2. PGS. TS. Lê Thị Thu Hiền

Hà Nội - 2022

LỜI CAM ĐOAN

Tôi xin cam đoan đề tài nghiên cứu trong luận văn này là công trình nghiên cứu của tôi và nhóm nghiên cứu dựa trên những tài liệu, số liệu do chính tôi và các thành viên trong nhóm tự tìm hiểu và nghiên cứu. Chính vì vậy, các kết quả nghiên cứu đảm bảo trung thực và khách quan nhất. Đồng thời, kết quả này chưa từng xuất hiện trong bất kỳ một nghiên cứu nào của các nhóm nghiên cứu khác. Các số liệu, kết quả nêu trong luận văn là trung thực, nếu sai tôi hoàn chịu trách nhiệm.

Tác giả luận văn

LỜI CẢM ƠN

Lời đầu tiên, em xin gửi lời cảm ơn chân thành nhất đến toàn bộ Quý Thầy Cô của Học viện Khoa học và Công nghệ vì đã truyền tải những kiến thức quý báu, những kinh nghiệm đầy thực tế cho em. Mặc dù đại dịch COVID-19 gây ra nhiều khó khăn cho quá trình học tập, sự tâm huyết và những truyền đạt tận tâm từ Quý Thầy Cô đã giúp em rất nhiều trong suốt hai năm học vừa qua.

Với tất cả sự trân trọng và quý mến, em xin gửi lời cảm ơn sâu sắc tới PGS. TS. Lê Thị Thu Hiền, Phó Viện trưởng Viện Nghiên cứu hệ gen, Viện Hàn lâm Khoa học và Công nghệ Việt Nam. Cảm ơn cô đã hỗ trợ em vào những thời điểm khó khăn nhất, ở bên em ngay cả những lúc em đã định không tiếp tục theo đuổi việc học. Em cảm thấy vô cùng may mắn khi có được sự giúp đỡ của cô để có thể hoàn thành công việc nghiên cứu.

Em cũng vô cùng biết ơn sự giúp đỡ của TS. Nguyễn Tường Vân, cán bộ nghiên cứu Viện Công nghệ sinh học, Viện Hàn lâm Khoa học và Công nghệ Việt Nam. Để có thể hoàn thành luận văn này, sự hỗ trợ và giúp đỡ của cô Vân đặc biệt có ý nghĩa. Em thực sự cảm ơn cô đã động viên, giúp em tiếp tục hoàn thiện luận văn này.

Ngoài ra, em cũng xin được gửi lời cảm ơn tới các thành viên của Phòng Đa dạng sinh học hệ gen, Viện Nghiên cứu hệ gen, đặc biệt cảm ơn ThS. Lưu Hàn Ly, ThS. Phạm Lê Bích Hằng, CN. Vũ Thị Trinh và CN. Nguyễn Thị Bích Ngọc, vì đã hỗ trợ em rất nhiều trong quá trình nghiên cứu và học hỏi, để em có thể hoàn thiện luận văn này; cảm ơn ThS. Nguyễn Nhật Linh vì những hướng dẫn quý giá vào những ngày đầu em tiếp cận với hướng nghiên cứu tin sinh học; cảm ơn TS. Vũ Tuấn Nam vì những hỗ trợ trong quá trình nghiên cứu, cũng như những kinh nghiệm quý báu trong cuộc sống.

Nghiên cứu này được thực hiện với sự hỗ trợ từ đề tài “Giải trình tự và phân tích hệ phiên mã (transcriptome) ở sâm Ngọc Linh (*Panax vietnamensis* Ha et Grushv.).

Sự giúp đỡ, động viên và chia sẻ của gia đình, các thầy cô và bạn bè hỗ trợ em rất nhiều trong quá trình hoàn thành luận văn này. Em vô cùng trân trọng và biết ơn tất cả mọi người!

Hà Nội, ngày 20 tháng 12 năm 2022
Tác giả luận văn

MỤC LỤC

DANH MỤC CÁC TỪ VIẾT TẮT	vi
DANH MỤC HÌNH	vii
DANH MỤC BẢNG	viii
MỞ ĐẦU	9
CHƯƠNG 1. TỔNG QUAN TÀI LIỆU	11
1.1 GIỚI THIỆU VỀ CHI <i>PANAX</i> VÀ SÂM NGỌC LINH	11
1.2 GIẢI TRÌNH TỰ DNA THỂ HỆ MỚI TRÊN THỂ GIỚI	16
1.3 NGHIÊN CỨU GIẢI TRÌNH TỰ HỆ PHIÊN MÃ Ở CÁC LOÀI THUỘC CHI <i>PANAX</i>	17
1.4 VAI TRÒ CỦA CÁC CÔNG CỤ TIN SINH TRONG PHÂN TÍCH TRÌNH TỰ HỆ PHIÊN MÃ	18
CHƯƠNG 2. VẬT LIỆU VÀ PHƯƠNG PHÁP NGHIÊN CỨU	21
2.1 VẬT LIỆU.....	21
2.1.1 Các mẫu sâm nghiên cứu.....	21
2.1.2 Hóa chất, thiết bị.....	21
2.2 PHƯƠNG PHÁP NGHIÊN CỨU.....	22
2.2.1 Thu mẫu.....	22
2.2.1.1 Phương pháp thu mẫu và bảo quản mẫu nhằm giải trình tự hệ phiên mã .	22
2.2.1.2 Phương pháp tách chiết RNA tổng số	22
2.2.1.3 Phương pháp tổng hợp cDNA và chuẩn bị thư viện cDNA	22
2.2.1.4 Giải trình tự hệ phiên mã.....	22
2.2.2 Phân tích và lắp ráp trình tự các hệ phiên mã	23
2.2.2.1 Kiểm tra chất lượng dữ liệu thô sau khi giải trình tự và tiền xử lý số liệu	23
2.2.2.2 Lắp ráp <i>de novo</i> các đoạn đọc	23
2.2.2.3 Phân nhóm các đoạn trình tự thành các unigene (Clustering).....	23
2.2.2.4 Dự đoán khung đọc mở (ORF).....	24
2.2.2.5 Ước lượng độ phong phú.....	24
2.2.3 Chú giải chức năng hệ phiên mã	24
2.2.3.1 Chú giải hệ phiên mã dựa trên cơ sở dữ liệu GO	24

2.2.3.2	Chú giải dựa trên cơ sở dữ liệu EggNOG	25
2.2.3.3	Chú giải dựa trên cơ sở dữ liệu NT và NR của NCBI.....	25
2.2.3.4	Chú giải dựa trên cơ sở dữ liệu KEGG	25
2.2.3.5	Chú giải dựa trên cơ sở dữ liệu UniProt.....	26
2.2.3.6	Chú giải dựa trên cơ sở dữ liệu Pfam	26
CHƯƠNG 3. KẾT QUẢ VÀ THẢO LUẬN		27
3.1 GIẢI TRÌNH TỰ HỆ PHIÊN MÃ ĐẶC HIỆU MÔ LÁ VÀ THÂN RỄ SÂM NGỌC LINH 4 NĂM TUỔI		27
3.1.1	Kết quả tách chiết, tinh sạch RNA tổng số và xây dựng các thư viện cDNA.....	27
3.1.2	Kết quả giải trình tự hệ phiên mã của mô lá và mô thân rễ sâm Ngọc Linh 4 năm tuổi	28
3.2 PHÂN TÍCH VÀ LẮP RÁP CÁC HỆ PHIÊN MÃ.....		29
3.2.1	Kết quả kiểm tra chất lượng các đoạn đọc sau khi trimming	29
3.2.2	Kết quả lắp ráp <i>de novo</i> các hệ phiên mã	31
3.2.3	Kết quả phân nhóm các đoạn trình tự thành các unigene.....	33
3.2.4	Kết quả dự đoán khung đọc mở.....	34
3.2.5	Kết quả ước lượng độ phong phú	35
3.3 CHÚ GIẢI HỆ PHIÊN MÃ CỦA CÁC MÔ SÂM NGỌC LINH 4 NĂM TUỔI		36
3.3.1	Chú giải hệ phiên mã mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi	36
3.3.1.1	Kết quả chú giải dựa trên cơ sở dữ liệu GO	36
3.3.1.2	Kết quả chú giải dựa trên cơ sở dữ liệu EggNOG.....	37
3.3.1.3	Kết quả chú giải dựa trên cơ sở dữ liệu NT và NR của NCBI.....	37
3.3.1.4	Kết quả chú giải dựa trên cơ sở dữ liệu KEGG.....	38
3.3.1.5	Kết quả chú giải dựa trên cơ sở dữ liệu UniProt	38
3.3.1.6	Kết quả chú giải dựa trên cơ sở dữ liệu Pfam	38
3.3.2	Kết quả chú giải các unigene của mẫu mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi.....	47
3.3.3	Tổng hợp, phân tích và so sánh dữ liệu hệ phiên mã ở mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi.....	48
CHƯƠNG 4. KẾT LUẬN VÀ KIẾN NGHỊ.....		55
4.1	KẾT LUẬN	55

4.2	KIẾN NGHỊ.....	55
	TÀI LIỆU THAM KHẢO	57

DANH MỤC CÁC TỪ VIẾT TẮT

Chữ viết tắt	Tên đầy đủ (tiếng Anh)	Tên đầy đủ (tiếng Việt)
BLAST	Basic local alignment search tool	Công cụ tìm kiếm và so sánh các trình tự tương đồng
Bp	Base pair	Cặp base
DNA	Deoxyribonucleic acid	Acid deoxyribonucleic
Đtg	<i>et al.</i>	Đồng tác giả
EST	Expressed sequence tag	Các đoạn trình tự gen biểu hiện
FastQC	Fast quality check	Công cụ kiểm tra đánh giá chất lượng dữ liệu
IUCN	International Union for Conservation of Nature and Natural Resources	Liên minh Quốc tế Bảo tồn thiên nhiên và Tài nguyên thiên nhiên
Kb	Kilo base	Kilo base (1.000 bp)
NCBI	National Center for Biotechnology Information	Trung tâm Tin sinh học Quốc gia
RNA-Seq	Whole transcriptome shotgun sequencing	Giải trình tự toàn bộ hệ phiên mã
ORF	Open reading frame	Khung đọc mở
SSR	Simple sequence repeat	Trình tự lặp lại đơn giản

DANH MỤC HÌNH

Hình 2.1 Hình ảnh một số mẫu sâm Ngọc Linh thu thập	21
Hình 2.2 Một số bước phân tích, lắp ráp và chú giải trình tự hệ phiên mã	23
Hình 3.1 Kiểm tra kết quả tách chiết và tinh sạch RNA tổng số trên gel agarose	27
Hình 3.2 So sánh dữ liệu thô và dữ liệu sau trimming ở mẫu mô lá (L4.1) và thân rễ (C4.1) của sâm Ngọc Linh 4 năm tuổi	30
Hình 3.3. Kết quả chú giải và phân nhóm các gen thuộc hệ phiên mã của mẫu mô lá (L4.1) dựa trên cơ sở dữ liệu GO	40
Hình 3.4. Kết quả chú giải và phân nhóm các gen thuộc hệ phiên mã của mẫu mô lá (L4.1) dựa trên cơ sở dữ liệu EggNOG	41
Hình 3.5 Tỷ lệ các unigene của hệ phiên mã mẫu mô lá (L4.1) và thân rễ (C4.1) được chú giải trên các cơ sở dữ liệu	48
Hình 3.6 Chú giải GO cho unigene sâm Ngọc Linh 4 năm tuổi	49
Hình 3.7 Chú giải KEGG của hệ phiên mã sâm Ngọc Linh 4 năm tuổi	50
Hình 3.8 Biểu đồ venn biểu diễn số unigene được chú giải bởi các cơ sở dữ liệu của hệ phiên mã mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi.....	51

DANH MỤC BẢNG

Bảng 1.1 Danh sách các loài trong chi <i>Panax</i> trên thế giới	11
Bảng 1.2 Một số đặc điểm hình thái của các thứ/loài thuộc chi <i>Panax</i> ở nước ta sử dụng trong định loại hình thái các mẫu thu thập	13
Bảng 3.1 Nồng độ và độ sạch (A260/A280) của một số mẫu RNA sâm Ngọc Linh sau tách chiết và tinh sạch.....	28
Bảng 3.2 Thống kê các bộ dữ liệu thô khi giải trình tự các thư viện cDNA.....	29
Bảng 3.3 Thống kê các bộ dữ liệu thu được sau khi trimming	29
Bảng 3.4 Kết quả thống kê của contig được lắp ráp đầu tiên.....	32
Bảng 3.5 Kết quả thống kê quá trình lắp ráp các đoạn đọc ở mẫu mô lá (L4.1) và thân rễ (C4.1).....	32
Bảng 3.6 Kết quả thống kê của contig unigene	34
Bảng 3.7 Kết quả thống kê dự đoán ORF của các mẫu mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi	35
Bảng 3.8 Tỷ lệ lắp ráp các đoạn đọc của mẫu mô lá và thân rễ sâm.....	36
Bảng 3.9 Kết quả chú giải và phân nhóm các gen thuộc hệ phiên mã của mẫu mô lá (L4.1) và thân rễ (C4.1) sâm Ngọc Linh 4 năm tuổi dựa trên cơ sở dữ liệu NT/NR của NCBI.....	42
Bảng 3.10 Kết quả chú giải và phân nhóm các gen thuộc hệ phiên mã của mẫu mô lá (L4.1) và thân rễ (L4.1) sâm Ngọc Linh 4 năm tuổi dựa trên cơ sở dữ liệu KEGG của NCBI.....	44
Bảng 3.11 Kết quả chú giải và phân nhóm các gen thuộc hệ phiên mã của mẫu mô lá (L4.1) và thân rễ (L4.1) sâm Ngọc Linh 4 năm tuổi dựa trên cơ sở dữ liệu UniProt của NCBI.....	45
Bảng 3.12 Kết quả chú giải và phân nhóm các gen thuộc hệ phiên mã của mẫu mô lá (L4.1) và thân rễ (L4.1) sâm Ngọc Linh 4 năm tuổi dựa trên cơ sở dữ liệu Pfam của NCBI.....	46
Bảng 3.13 Tỷ lệ các unigene được chú giải trong hệ phiên mã các mẫu mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi	47
Bảng 3.14 Khả năng chú giải các unigene của hệ phiên mã sâm Ngọc Linh 4 năm tuổi trên các cơ sở dữ liệu khác nhau.....	48
Bảng 3.15 Danh sách transcript chiếm ưu thế trong hệ phiên mã của mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi	52

MỞ ĐẦU

Với nền y học cổ truyền trên thế giới nói chung và ở Việt Nam nói riêng, từ lâu sâm đã là nguồn dược liệu quý giá. Ngày nay, tên gọi sâm được dùng phổ biến để chỉ một số loài thực vật, trong đó chủ yếu là các loài thuộc chi *Panax* L. và *Eleutherococcus* Maxim thuộc họ Araliaceae. Giá trị kinh tế của sâm về mặt thương mại được ước tính vượt 2,1 tỷ đô la Mỹ trên quy mô toàn cầu [1]. Các loài sâm được sử dụng phổ biến nhất trong chi này là sâm Hàn Quốc (*Panax ginseng*), sâm Mỹ (*Panax quinquefolius*), sâm Trung Quốc (*Panax notoginseng*), sâm Nhật Bản (*Panax japonicus*), sâm Himalaya (*Panax pseudoginseng*) và sâm Việt Nam (*Panax vietnamensis* Ha et Grushv.) [2]. Nhiều loài sâm đang bị đe dọa trong tự nhiên, do đó các sáng kiến phục vụ phát triển hiệu quả và bền vững là đặc biệt cần thiết để bảo tồn đa dạng sinh học ở cả cấp độ loài và cấp độ hệ gen, đồng thời đáp ứng nhu cầu về sản phẩm dược liệu. *P. vietnamensis*, hay còn được biết đến với tên gọi sâm Ngọc Linh, là loài dược liệu đặc hữu ở nước ta. Sâm Ngọc Linh có giá trị cao về khoa học và kinh tế, tuy nhiên các nghiên cứu còn hạn chế. Cũng giống như các loài thuộc chi *Panax*, sâm Ngọc Linh mang các hợp chất triterpene saponin tự nhiên, hay còn được biết đến với tên gọi là ginsenoside. Những hợp chất này có hoạt tính tốt, đồng thời mang tác dụng chống suy nhược cơ thể, oxy hóa, đái tháo đường, ung thư, tăng cường thể lực, bảo vệ hệ thần kinh. Do việc khai thác quá mức cũng như phân bố tương đối hẹp, sâm Ngọc Linh trở nên khó phát triển trong môi trường tự nhiên và đã được thêm vào Danh lục đỏ của Liên minh Quốc tế Bảo tồn thiên nhiên và Tài nguyên thiên nhiên (IUCN) và Sách đỏ Việt Nam (2007) [3]. Do đó, những nghiên cứu di truyền cũng như đặc tính hóa học và dược tính của sâm Ngọc Linh đã và đang được thực hiện, góp phần phát triển loài thực vật quý hiếm này ở nước ta.

Với mục đích bổ sung thêm những thông tin hữu ích về hệ phiên mã sâm Ngọc Linh, qua đó hỗ trợ đánh giá nguồn gen, giúp bảo tồn, khai thác và sử

dụng bền vững sâm Ngọc Linh ở Việt Nam, nghiên cứu “**Ứng dụng tin sinh học trong phân tích hệ phiên mã cây sâm Ngọc Linh (*Panax vietnamensis* Ha et Grushv.) 4 năm tuổi**” đã được thực hiện với mục tiêu và những nội dung cụ thể như sau:

1 MỤC TIÊU

Nghiên cứu đặt mục tiêu phân tích được kết quả giải trình tự và xây dựng được cơ sở dữ liệu của hệ phiên mã đặc hiệu mô (lá, thân rễ) ở sâm Ngọc Linh 4 năm tuổi.

2 NỘI DUNG

- Giải trình tự hệ phiên mã đặc hiệu mô thân rễ và lá của sâm Ngọc Linh 4 năm tuổi;
- Phân tích, lắp ráp, hiệu chỉnh trình tự các hệ phiên mã mô lá và thân rễ;
- Chú giải hệ phiên mã mô thân rễ và lá.

CHƯƠNG 1. TỔNG QUAN TÀI LIỆU

1.1 GIỚI THIỆU VỀ CHI *PANAX* VÀ SÂM NGỌC LINH

Chi *Panax* L., hay chi Sâm thuộc họ Araliaceae (Ngũ gia bì/ Nhân sâm), bộ Apiales (Hoa tán), là một trong những chi bao gồm nhiều loài dược liệu quý. *Panax* phân bố tự nhiên chủ yếu ở một số quốc gia thuộc khu vực Bắc Mỹ, châu Á, đặc biệt là vùng Đông Á. Nhiều loài *Panax* chứa các ginsenoside, hay còn được biết tới là các hợp chất tự nhiên triterpene saponin. Các hợp chất này có cấu tạo phân tử phức tạp; có lợi cho sức khỏe con người: bảo vệ hệ thần kinh, chống đái tháo đường, oxy hóa, ung thư và tăng cường thể lực [3, 4]. Do vùng phân bố hạn chế và giá trị kinh tế cùng nhu cầu rất cao nên nhiều loài thuộc chi *Panax* đã bị khai thác quá nhiều và khan hiếm trong môi trường tự nhiên. Nhiều quốc gia cũng đã triển khai các chương trình khai thác, bảo tồn cũng như sử dụng nguồn gen bền vững đặc biệt quý này, chẳng hạn như sâm Trung Quốc, Hàn Quốc và Mỹ được trồng thương mại và phổ biến nhất [5]. Gần đây, 12 loài và 2 thứ thuộc chi này đã được mô tả (Bảng 1.1) [5].

Bảng 1.1 Danh sách các loài trong chi *Panax* trên thế giới

Tên khoa học	Khu vực địa lý
<i>P. bipinnatifidus</i> Seem.	
<i>P. bipinnatifidus</i> var. <i>angustifolius</i>	
<i>P. bipinnatifidus</i> var. <i>bipinnatifidus</i>	
<i>P. ginseng</i> C.A.Mey.	Hàn Quốc (Nhân sâm)
<i>P. japonicus</i> (T.Nees) C.A.Mey.	Nhật Bản
<i>P. notoginseng</i> (Burkill) F.H.Chen	Trung Quốc
<i>P. pseudoginseng</i> Wall.	Himalaya

Tên khoa học	Khu vực địa lý
<i>P. quinquefolius</i> L.	Mỹ
<i>P. sokpayensis</i> Shiva K.Sharma et Pandit	
<i>P. stipuleanatus</i> H.T.Tsai et K.M.Feng	
<i>P. trifolium</i> L.	
<i>P. vietnamensis</i> Ha et Grushv.	Việt Nam
<i>P. wangianus</i> S.C.Sun	
<i>P. zingiberensis</i> C.Y.Wu et Feng	

Ở Việt Nam, các nghiên cứu về xác định thành phần và phân loại loài cho thấy, chi *Panax* gồm 3 loài mọc trong tự nhiên và có thể sử dụng trong dược liệu, bao gồm *P. vietnamensis* Ha et Grushv. (sâm Ngọc Linh), *P. bipinnatifidus* Seem. (sâm Vũ diệp) và *P. stipuleanatus* Tsai et Feng (Tam thất hoang) [6, 7].

Dù được tìm thấy lần đầu tiên vào năm 1973, phải tới năm 1985, Sâm Ngọc Linh mới được chính thức công nhận là loài thực vật mới [7]. Cây sâm này có các tên gọi khác nhau như sâm khu Năm (K5), sâm Việt Nam, trúc tiết sâm (hay sâm đốt trúc), củ Ngải rơm hay cây Thuốc giầu. Sâm Ngọc Linh mọc dày thành đám dưới tán rừng dọc theo các suối ẩm trên đất nhiều mùn, tập trung ở khu vực Ngọc Linh (thuộc tỉnh Kon Tum và Quảng Nam). Đây là dãy núi cao thứ hai ở nước ta (đỉnh cao nhất cao 2.598 m). Hai thứ khác của *P. vietnamensis* được phát hiện gần đây là *P. vietnamensis* var. *fuscidiscus* K.Komatsu, S.Zhu et S.Q.Cai (sâm Lai Châu) phân bố ở Lai Châu, *P. vietnamensis* Ha et Grushvitzky var. *langbianensis* N.V.Duy, V.T.Tran & L.N.Trieu (sâm Lang Bian) được phát hiện ở vùng núi Langbian, Lâm Đồng. Ngoài ra, *Panax* sp. (*Panax* sp. Puxailaileng) được phát hiện phân bố ở vùng núi Puxailaileng, tỉnh Nghệ An và được công nhận là một thứ của Ngọc Linh, tuy nhiên cần các nghiên cứu sâu

hơn [9]. Tương tự với nhiều loài thuộc chi *Panax*, sâm Ngọc Linh rất hiếm khi được tìm thấy trong môi trường tự nhiên do khai thác quá mức cũng như phân bố hạn chế, có mặt trong Sách đỏ Việt Nam [6].

Bảng 1. mô tả một số đặc điểm hình thái của các thứ/loài thuộc chi *Panax* ở nước ta (bao gồm sâm Ngọc Linh, sâm Vũ diệp, sâm Nghệ An, sâm Lai Châu, Tam thất hoang) [6, 8, 9].

Bảng 1.2 Một số đặc điểm hình thái của các thứ/loài thuộc chi *Panax* ở nước ta sử dụng trong định loại hình thái các mẫu thu thập

Đặc điểm	Sâm Ngọc Linh	Sâm Lai Châu	Sâm Nghệ An	Sâm Vũ diệp	Tam thất hoang
Chiều cao cây (cm)	30-110	40-80		30-100	25-75
Rễ	Có các đốt và có thể phân nhánh, nằm ngang (đường kính khoảng 1 tới 2 cm)	Thân rễ mập, nạc, hợp trục, nằm hơi chéch hoặc ngang, không rễ phụ (đường kính: khoảng 1,5-2,2 cm)	Thân rễ mọc bò ngang, mang nhiều rễ nhánh và củ, (đường kính: 1-2 cm)	Thân rễ lớn, phương ngang, thường ở trên mặt đất, (đường kính khoảng 1,5 tới 3,5 cm)	Rễ mập, nằm ngang, thường ít phân nhánh, một số chỗ lõm do vết thân, (đường kính khoảng 1,5 tới 3 cm)
Thân (mang lá)	Phụ thuộc số lượng nhánh ở rễ, phần mang lá có thể có từ 1 tới 5 thân	Thân đơn độc, nhẵn, mọc thẳng đứng, vỏ màu lục, cao 0,3-0,7 m, phần giữa xốp khi tươi và rỗng khi khô	Thân kí sinh, thẳng đứng, nhỏ, có đường kính thân 4-8 mm, màu xanh lục hoặc hơi tím,	Phụ thuộc số lượng nhánh của rễ, phần mang lá có từ 1 tới 3 thân, đường kính thân: 0,3-0,6 cm	Thường có 1 thân, ít khi 2 hoặc 4 thân. Mọc thẳng, nhẵn, đường kính từ 0,3 tới 0,6 cm

Đặc điểm	Sâm Ngọc Linh	Sâm Lai Châu	Sâm Nghệ An	Sâm Vũ diệp	Tam thất hoang
Lá	Dạng kép hình chân vịt, mọc vòng (3 tới 5 lá), dạng bầu dục, nhọn và thuôn 2 đầu, kích thước khoảng 6 tới 14 x 2,3 tới 4 cm, mép có thêm răng cưa	Lá mọc vòng 3 tới 6 lá ở đỉnh thân, dạng kép hình chân vịt, mang 5 tới 7 lá chét mỏng, dạng màng. Cuống lá dài khoảng 7 tới 14 cm.	Lá kép chân vịt mọc vòng với 3 tới 5 nhánh lá. Cuống lá kép: 6-12 mm, 5 lá, lá ở chính giữa (12-15 cm). Lá hình bầu dục, mép răng cưa, chóp nhọn, lông gai ở hai mặt lá	Lá kép chân vịt, thường gồm 3 cái mọc vòng ở ngọn. Lá chét (từ 3 tới 5), có răng cưa	Lá kép hình chân vịt (khoảng 1 tới 3 lá), mọc vòng ở ngọn, cuống khoảng 5 tới 10 cm. Lá cuống khá ngắn, phiến hình hoặc mác thuôn, nhọn 2 đầu, kích thước 5x2 -13x4 cm, mép và cửa thùy có răng cưa, lông ở gân lá phía trên
Cụm hoa	Tán đơn hoặc kép (có thể có 1 - 2 tán phụ), mọc ở ngọn, chiều dài của cuống 15-30 cm, cao vượt tán do dài hơn cuống lá	Mọc đơn ở giữa vòng và đỉnh thân. Cuống dài khoảng 25 cm. Cụm hoa là tán, gần hình cầu, đường kính đến từ 2,5 tới 4 cm, có 70 tới 100 hoa (hoặc hơn)	Tán đơn, dưới các lá thẳng với thân, cuống tán hoa 10 tới 20 cm (có thể kèm 1 tới 4 tán phụ hay hoặc 1 hoa riêng lẻ). Có 60-100 hoa trong một tán.	Tán mọc ngọn, đơn, cuống 5-510 cm, cụm hoa có 20 tới 90 hoa, cuống mảnh, dài khoảng 1 tới 1,5 cm	Tán đơn, mọc ngọn (ít khi có tán phụ nhỏ). Cuống cụm cao tương đương hoặc hơn tán lá, dài 5 tới 10 cm. Hoa trên một tán: 30-80
Hoa	Cuống ngắn, trắng xanh, với 5 đài hoa, 5 nhị và 5	Hoa vàng nhạt hoặc lục, rộng 3-4 mm. Bầu ở dưới, gồm 2 lá noãn,	Cuống hoa ngắn 1 tới 1,5 cm, lá đài 5, cánh hoa 5, màu vàng nhạt,	Màu vàng xanh, gồm 5 nhị, 5 cánh, 5 lá đài nhỏ. Đầu vòi	Hoavàng xanh, 5 lá đài khá nhỏ, có 5 cánh, 5 nhị hoa, bầu có 2 ô, đầu nhụy

Đặc điểm	Sâm Ngọc Linh	Sâm Lai Châu	Sâm Nghệ An	Sâm Vũ diệp	Tam thất hoang
	cánh. Bầu hai ô, vòi nhụy chẻ chia đôi ở đầu	hợp với vòi nhụy hình thành bầu, chẻ đôi	nhị 5, bầu 1 ô với 1 vòi nhụy	nhụy chẻ chia đôi. Bầu có 2 ô.	chẻ làm đôi.
Quả	Quả hình cầu, mọng (đường kính khoảng 0,5 tới 0,6 cm). Chín đỏ và thường xuất hiện chấm đen ở đỉnh.		Tập trung ở trung tâm của tán, dài 0,8-1cm, rộng 0,5-0,6 cm.	Quả hình cầu hoặc cầu dẹt, đường kính 0,6-1,2 cm, có màu đỏ khi chín	Hình cầu dẹt, mọng, đường kính khoảng 0,6 tới 1,2 cm và thời điểm chín có màu đỏ
Hạt	1 - 2 hạt, hạt nhỏ, hơi tròn hoặc gần giống hình thận. Không nhăn			Hạt 2, hơi tròn và có màu xám trắng, vỏ khá cứng và có rốn hạt	1 - 2 hạt, khá giống với hạt đậu, màu xám trắng, vỏ cứng và có rốn hạt
Mùa hoa, quả	Mùa hoa vào tháng 4-5, quả 6-9			Hoa tháng 4-5 và quả tháng 5-9 (10)	Hoa tháng 4-5. Quả tháng 5-9 (10)

Với giá trị khoa học và kinh tế rất cao, trong thời gian qua, Kon Tum và Quảng Nam đã tích cực phát triển và bảo tồn sâm Ngọc Linh. Việc trồng sâm Ngọc Linh trên vùng bán sinh trưởng hoặc sinh trưởng tự nhiên cho kết quả cây sâm phát triển và sinh trưởng tương đối tốt. Từng hộ gia đình được giao sâm giống để quản lý và chăm sóc. Nếu được tập trung phát triển, cách làm này sẽ

mang lại hiệu quả thiết thực và góp phần xóa đói giảm nghèo cho các hộ đồng bào dân tộc thiểu số [10].

1.2 GIẢI TRÌNH TỰ DNA THỂ HỆ MỚI TRÊN THẾ GIỚI

Sự thành công trong nghiên cứu giải trình tự hệ gen của con người đã mang đến một thời kỳ phát triển mạnh mẽ đối với nghiên cứu khoa học. Các công nghệ giải trình tự mới, phức tạp và hiệu quả hơn lần lượt ra đời, ứng dụng cho nghiên cứu và phát triển khoa học. Trong đó, một công nghệ có sức ảnh hưởng mạnh mẽ ở quy mô toàn cầu chính là giải trình tự gen thế hệ mới (NGS - next generation sequencing). Công nghệ này cho phép giải trình tự nhanh chóng và chính xác, được sử dụng cho nhiều ứng dụng như giải trình tự toàn bộ hệ gen sinh vật (whole genome sequencing), hệ phiên mã transcriptome (RNA-seq) hay hệ gen biểu hiện (whole exome sequencing). Tiếp nối từ Dự án giải trình tự hệ gen người, vào năm 2005, 454 sequencing đã ra mắt thị trường. Năm 2006, Genome Analyzer được cho ra đời bởi Solexa và SOLiD được phát triển bởi Agencourt. Các hệ thống này có các ưu điểm là độ chính xác rất cao, giá thành hạ thấp, đồng thời có dung lượng rất lớn so với trước đó. Các tổ chức đó sau này được mua lại bởi Applied Biosystems/ Agencourt vào năm 2006), Roche/ 454 và Illumina/ Solexa vào năm 2007. Đến nay, nhiều hệ thống máy NGS đã được phát triển, có thể kể tới Roche/ 454; Illumina NextSeq, MiSeq, NovaSeq, HiSeq.; Pacific Biosciences/ RS; Applied Biosystem/ SOLiD; Life technologies/ Ion Torrent PGM, Life technologies/ Ion Proton... nhờ đó giúp giải trình tự cho toàn bộ hệ gen một cách chính xác và nhanh chóng [11, 12, 13, 14, 15].

Các hệ thống giải trình tự thế hệ mới trên hoạt động dựa vào nguyên tắc gắn nối hoặc tổng hợp. Công nghệ NGS được cải tiến trên cơ sở các bước như sau: (i) chuẩn bị mẫu đầu vào, (ii) giải trình tự, (iii) lắp ráp hệ gen, (iv) chú giải và so sánh. Công đoạn giải trình tự bao gồm các bước như sau: (i) chuẩn bị và gắn lên giá bám các đoạn DNA; (ii) khuếch đại các DNA trên nhờ môi đặc hiệu; (iii) giải trình tự bằng gắn nối hoặc tổng hợp [16, 17]. Với ưu thế đọc nhanh

chóng, trình tự đọc được xử lý với dung lượng rất lớn và hiệu quả cao, các hệ thống NGS ngày càng được ứng dụng rộng rãi [18, 19, 20].

1.3 NGHIÊN CỨU GIẢI TRÌNH TỰ HỆ PHIÊN MÃ Ở CÁC LOÀI THUỘC CHI *PANAX*

Bên cạnh sự phát triển mạnh mẽ của NGS, việc nghiên cứu giải trình tự hệ phiên mã cho các loài ở chi *Panax* đã và đang được tiến hành trong thời gian gần đây. Thống kê cho thấy, các nghiên cứu hiện tại chủ yếu tập trung ở một số khu vực, quốc gia có phân bố cây sâm và thường xuyên được sử dụng làm dược phẩm, chẳng hạn như Hàn Quốc hay Trung Quốc.

Trong năm 2010, tại Viện Phát triển thảo dược Bắc Kinh, Sun và đtg đã công bố nghiên cứu đầu tiên về giải trình tự hệ phiên mã của *P. quinquefolius*. Nhóm tác giả sử dụng hệ thống GS FLX Titanium, sử dụng công nghệ 454 pyrosequencing, qua đó thu được khoảng 200.000 đoạn đọc chất lượng tốt, có độ dài mỗi đoạn đọc trung bình của thư viện cDNA rễ lên tới 427 bp. Kết quả, trên 30.000 unigene đã được xác định và khoảng 2/3 số trình tự này có thể tìm thấy từ các ngân hàng trình tự [21].

Năm 2011, Luo và đtg đã công bố hệ phiên mã của *P. notoginseng* thuộc chi *Panax* [22]. Các nhân tố điều hòa dịch mã (WRKY, Myb, bHLH, homeobox) và các trình tự SSR ở *P. notoginseng* cũng được nghiên cứu từ dữ liệu hệ phiên mã.

Dựa trên số liệu từ 2 lần chạy hệ thống 454 pyrosequencing, vào năm 2013, Li và đtg đã thu được 45.849, 6.172, 4.041 và 3.273 trình tự mã hóa từ thư viện cDNA lần lượt của rễ, thân, lá và hoa của *P. ginseng* [23]. Nghiên cứu của Cao và đtg (2015) đã phát hiện được 2 unigene mã hóa geranyl diphosphate synthase (GPS) ở rễ sâm *P. ginseng* 4 năm tuổi [24]. Dựa trên các trình tự hệ gen tham chiếu của *P. ginseng* giống “Chunpoong” đã được công bố, các bộ dữ liệu hệ phiên mã và chú giải chức năng, Jayakodi và đtg đã xây dựng cơ sở dữ liệu về hệ gen của *P. ginseng* tại <http://ginsengdb.snu.ac.kr/> [25]. Đây là nền

tăng truy cập mở đầu tiên cung cấp dữ liệu toàn diện của *P. ginseng*, là nguồn tài nguyên quý giá cho nhiều lĩnh vực nghiên cứu liên quan đến *P. ginseng* và các loài khác thuộc bộ Apiales cũng như cho cộng đồng các nhà khoa học nghiên cứu thực vật nói chung.

Trên đối tượng sâm Nhật Bản *P. japonicus*, Rai và đtg (2016) đã tiến hành phân tích trình tự hệ phiên mã, đồng thời so sánh các gen tiềm năng liên quan đến sinh tổng hợp các saponin với các loài sâm khác [26]. Kết quả giải trình tự RNA bằng hệ thống Illumina thu được 135.235 unigene. Trên đối tượng *P. vietnamensis*, Vu và đtg (2020) đã công bố kết quả giải trình tự hệ phiên mã của *P. vietnamensis*, chú giải các unigene sử dụng 7 cơ sở dữ liệu bao gồm Clusters of Orthologous Groups (COG), GO, KEGG, Clusters of Orthologous Groups (KOG), Pfam, Swissprot và NR. Mẫu phân tích là mẫu mô hỗn hợp của các mẫu sâm không xác định tuổi và tỉ lệ chú giải đạt 35,49% với 31.686 unigene được phát hiện [27].

1.4 VAI TRÒ CỦA CÁC CÔNG CỤ TIN SINH TRONG PHÂN TÍCH TRÌNH TỰ HỆ PHIÊN MÃ

Các công nghệ giải trình tự thế hệ mới cho phép xác định dữ liệu lớn về gen/hệ gen trong thời gian ngắn. Kể từ khi các kỹ thuật giải trình tự hiện đại ra đời, việc xác định trình tự nucleic acid đã trở thành một công cụ phổ biến và thiết yếu trong mọi lĩnh vực khoa học sinh học. Tương ứng, lĩnh vực tin sinh học là trung tâm của việc giải thích và ứng dụng dữ liệu sinh học này. Sử dụng các phương pháp toán học và thống kê được triển khai bởi nhiều ngôn ngữ lập trình, các công cụ tin sinh học tổ chức, phân tích và giải thích thông tin sinh học ở cấp độ phân tử, tế bào và gen. Tin sinh học, sử dụng các công cụ tính toán, toán học và thống kê để thu thập, sắp xếp và phân tích dữ liệu trình tự gen lớn và phức tạp cũng như dữ liệu sinh học liên quan. Sức mạnh tổng hợp của NGS và tin sinh học rất quan trọng đối với phân tích, xử lý và đánh giá kết quả giải trình tự

hệ phiên mã, trong đó bao gồm ba bước chính: (1) kiểm soát chất lượng và chuẩn bị dữ liệu, (2) lắp ráp bộ gen; và (3) phân tích sau lắp ráp.

Các công nghệ giải trình tự hiện đại có thể tạo ra một số lượng lớn các lần đọc trình tự trong một thí nghiệm. Tuy nhiên, chưa có công nghệ giải trình tự nào là hoàn hảo và mỗi công cụ sẽ tạo ra các loại và số lượng lỗi khác nhau, ví dụ các nucleotide được xác định không chính xác. Tỷ lệ base được xác định sai này sẽ phụ thuộc vào giới hạn kỹ thuật của từng nền tảng giải trình tự. Việc xác định và loại trừ các loại lỗi có thể ảnh hưởng đến việc phân tích tiếp theo là rất cần thiết. Do đó, kiểm soát chất lượng trình tự là bước đầu tiên thiết yếu trong phân tích, giúp tiết kiệm thời gian và chất lượng của các bước sau này. Việc kiểm soát chất lượng cơ sở dữ liệu thô sau giải trình tự thường sử dụng công cụ FastQC [28] để đánh giá các tiêu chí như: chất lượng trình tự base, tỷ lệ adapter, điểm chất lượng trình tự, tỷ lệ base, tỷ lệ G-C, phân bố độ dài trình tự, tỷ lệ lặp lại của các trình tự, các trình tự biểu hiện quá mức và công cụ Trimmomatic [29] để loại bỏ các trình tự ngắn, ví dụ, dưới 36 base hoặc có điểm chất lượng thấp.

Sau khi có được bộ dữ liệu đã xử lý và làm sạch, các bước lắp ráp, phân đoạn trình tự thành unigene, dự đoán khung đọc mở và ước lượng độ phong phú đều yêu cầu bắt buộc các công cụ tin sinh tương ứng để giải quyết. Các công cụ hiện nay cho phép việc lắp ráp hệ gen *de novo*, giúp tạo các trình tự tham chiếu, ngay cả đối với hệ gen phức tạp hoặc đa bội, cung cấp thông tin hữu ích để lập bản đồ hệ gen của các sinh vật mới hoặc hoàn thiện hệ gen của các sinh vật đã biết, làm rõ các vùng rất giống nhau hoặc lặp đi lặp lại để lắp ráp *de novo* chính xác, cũng như xác định các biến thể cấu trúc và sắp xếp lại, chẳng hạn như xóa, đảo ngược hoặc chuyển vị.

Các công cụ tin sinh cũng đóng vai trò quan trọng trong quá trình chú giải các đoạn đọc đã lắp ráp. Chú giải hệ gen là quá trình lấy thông tin cấu trúc và chức năng của protein hoặc gen từ dữ liệu thô bằng cách sử dụng các phương pháp tin sinh để phân tích, so sánh, ước tính. Sau khi hệ gen được giải trình tự

và lắp ráp, các đoạn đọc cần được chú giải để có được những thông tin hợp lý hơn về các đặc điểm cấu trúc và vai trò chức năng của nó. Sử dụng các phương pháp chú giải gen, có thể dự đoán các gen hoặc protein thuộc hệ gen cụ thể. Chú giải chức năng của các gen hoặc protein mới này có thể được thực hiện bằng cách tìm kiếm sự giống nhau của chúng với các trình tự đã được xác định bằng thực nghiệm và có sẵn trong cơ sở dữ liệu.

Thông tin trình tự hệ gen được lưu trữ trong một số định dạng tệp như FASTA, GFF3 và GENBANK. Có nhiều định dạng tệp khác nhau để trình bày trình tự, cấu trúc và thông tin con đường liên quan đến gen và protein, đồng thời cơ sở để chọn và tải xuống một tệp cụ thể có sẵn trên cơ sở dữ liệu trực tuyến. Các cơ sở dữ liệu để thực hiện chú giải cũng đa dạng và phụ thuộc vào mục đích nghiên cứu. Một số cơ sở dữ liệu phổ biến bao gồm GenBank, EMBL; GO, KEGG, UniProtKB/Swiss-Prot, Pfam. Tương ứng với mỗi cơ sở dữ liệu, các công cụ tin sinh cũng được tối ưu thuật toán để khai thác dữ liệu trình tự gen/protein, đưa ra các thông tin chi tiết như kích thước, độ bao phủ, tỉ lệ tương đồng, số lượng gap, độ tin cậy, vị trí khung đọc mở... cũng như tham khảo các tài liệu nghiên cứu trước đó để làm tiền đề cho việc so sánh, đối chiếu.

CHƯƠNG 2. VẬT LIỆU VÀ PHƯƠNG PHÁP NGHIÊN CỨU

2.1 VẬT LIỆU

2.1.1 Các mẫu sâm nghiên cứu

Nghiên cứu sử dụng 6 mẫu sâm Ngọc Linh 4 năm tuổi (bao gồm 3 mẫu thân rễ C4.1, C4.2, C4.3 và 3 mẫu lá L4.1, L4.2, L4.3) thu thập tại Trại Sâm gốc Tắc Ngo, thuộc Trung tâm Sâm Ngọc Linh (Nam Trà My, Quảng Nam). Một số mẫu sâm Ngọc Linh 4 năm tuổi phục vụ phân tích hệ phiên mã được trình bày trên Hình 2.1.



Mẫu C4.1



Mẫu L4.1

Hình 2.1 Hình ảnh một số mẫu sâm Ngọc Linh thu thập

2.1.2 Hóa chất, thiết bị

Các hóa chất sử dụng được đảm bảo độ tinh khiết và chất lượng cho các thí nghiệm sinh học phân tử và có nguồn gốc từ các hãng như Thermo Fisher Scientific (Hoa Kỳ), Illumina (Hoa Kỳ), Qiagen (Hoa Kỳ)... Các máy móc, thiết bị chính được sử dụng thuộc Viện Hàn lâm Khoa học và Công nghệ Việt Nam và Bảo tàng Lịch sử Tự nhiên, Đại học Oslo, Na Uy. Ngoài ra, nghiên cứu sử

dụng các chương trình phần mềm chuyên dụng miễn phí hoặc được cấp phép tại Viện Nghiên cứu hệ gen thuộc Viện Hàn lâm Khoa học và Công nghệ Việt Nam và Bảo tàng Lịch sử Tự nhiên, Đại học Oslo, Na Uy.

2.2 PHƯƠNG PHÁP NGHIÊN CỨU

2.2.1 Thu mẫu

2.2.1.1 Phương pháp thu mẫu và bảo quản mẫu nhằm giải trình tự hệ phiên mã

Mẫu mô thân rễ và mô lá của cây sâm Ngọc Linh 4 năm tuổi được thu và bảo quản riêng biệt. Các cây được rửa sạch bằng nước cất đã khử trùng, thấm khô, cắt nhỏ và bảo quản trong dung dịch RNAlater tự pha. Mẫu mô trong RNAlater được giữ ở 4°C trong 24 giờ trước khi bảo quản thời gian dài trong tủ lạnh sâu -80°C.

2.2.1.2 Phương pháp tách chiết RNA tổng số

RNA tổng số của các mẫu lá, thân rễ sâm Ngọc Linh được tách chiết và tinh sạch sử dụng RNeasy Plant Mini Kit (Qiagen, Hoa Kỳ) theo hướng dẫn của nhà sản xuất.

RNA tổng số của các mẫu lá và thân rễ sâm Ngọc Linh được kiểm tra, đánh giá chất lượng bằng điện di trên gel agarose 0,8%. Để xác định nồng độ và độ sạch, các mẫu RNA được phân tích bằng máy quang phổ Biospectrometer (Eppendorf, Đức). RNA sau tách chiết và tinh sạch được bảo quản ở -80°C cho đến khi được sử dụng cho các thí nghiệm tiếp theo.

2.2.1.3 Phương pháp tổng hợp cDNA và chuẩn bị thư viện cDNA

Các đoạn mRNA sẽ được phiên mã ngược tạo thành các đoạn cDNA theo hướng dẫn của nhà sản xuất. Thư viện được chuẩn bị sử dụng TruSeq Stranded mRNA LT Sample Prep Kit (Illumina, Hoa Kỳ).

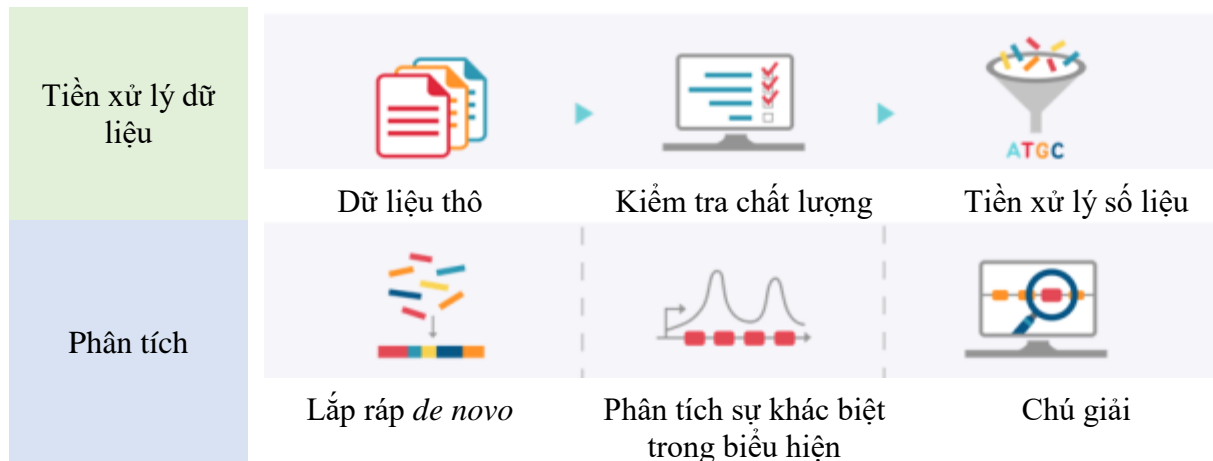
2.2.1.4 Giải trình tự hệ phiên mã

Các thư viện cDNA đạt yêu cầu chất lượng được sử dụng để tiến hành

giải trình tự hệ phiên mã bằng hệ thống NGS của Illumina.

2.2.2 Phân tích và lắp ráp trình tự các hệ phiên mã

Dữ liệu thô sau khi giải trình tự được kiểm tra chất lượng và tiền xử lý dữ liệu, sau đó trải qua các bước lắp ráp *de novo*, phân tích sự khác biệt trong biểu hiện và chú giải transcriptome (Hình 2.).



Hình 2.2 Một số bước phân tích, lắp ráp và chú giải trình tự hệ phiên mã

2.2.2.1 Kiểm tra chất lượng dữ liệu thô sau khi giải trình tự và tiền xử lý số liệu

Chất lượng các đoạn đọc thô sau khi giải trình tự được đánh giá lại, kiểm tra chất lượng bằng phần mềm FastQC v0.11.9. Phần mềm Trimmomatic v0.36 được sử dụng để tiền xử lý số liệu nhằm loại bỏ các trình tự là các adapter và các base chất lượng kém.

2.2.2.2 Lắp ráp *de novo* các đoạn đọc

Các đoạn đọc sau bước chọn lọc chất lượng sẽ được lắp ráp sử dụng phần mềm Trinity trinityrnaseq_r20140717 [30]. Các thông số được sử dụng trong lắp ráp *de novo* các đoạn đọc của mẫu phân tích đều được sử dụng theo mặc định của phần mềm. Trong quá trình lắp ráp, các đoạn đọc được chuẩn hóa *in silico* với độ bao phủ lớn nhất là 50.

2.2.2.3 Phân nhóm các đoạn trình tự thành các unigene (Clustering)

Đối với các gen đã được lắp ráp, contig dài nhất được lọc và chia thành các trình tự phiên mã sử dụng phần mềm CD-HIT-EST v4.6 [31, 32]. Các thông số được sử dụng trong phân tích đều là các thông số mặc định của phần mềm với ngưỡng giá trị phân nhóm là 0,9 tương đương độ tương đồng 90%, kích thước word là 8 base, sử dụng fast mode, loại bỏ các unigene với kích thước nhỏ hơn 200 bp... Các trình tự phiên mã sau lắp ráp được coi là các unigene.

2.2.2.4 Dự đoán khung đọc mở (ORF)

Các unigene đã thu được từ bước trên được sử dụng làm nguyên liệu đầu vào cho quá trình dự đoán khung đọc mở (Open Reading Frame - ORF). Phần mềm TransDecoder v3.0.1 [33] được sử dụng để tìm kiếm và phát hiện những vùng mã hóa trong trình tự hệ phiên mã. ORF được phát hiện bằng cách tìm kiếm các codon mở đầu và kết thúc. Các ORF được lựa chọn với kích thước trên 100 amino acid. Trong số đó, 500 ORF có kích thước lớn nhất được lựa chọn và một mô hình Markov bậc 5 đặc hiệu khung đọc được sử dụng. Các ORF sau đó sẽ được tính điểm và so sánh để tìm ra các vùng mã hóa. Ngoài ra, các ORF có kích thước tối thiểu lớn hơn 900 base cũng được lựa chọn.

2.2.2.5 Ước lượng độ phong phú

Bước ước lượng độ phong phú của các unigene được đánh giá sử dụng phần mềm RSEM v1.2.29 [34]. Các đoạn đọc sau khi lọc chất lượng được sử dụng để tính toán ước lượng dựa trên hệ phiên mã đã lắp ráp sử dụng phần mềm Bowtie tích hợp trong RSEM. Nhờ đó, RSEM có thể ước tính biểu hiện gen và tính toán độ phong phú của các unigene dựa trên số lượng các đoạn đọc. Các thông số được sử dụng trong ước lượng độ phong phú đều là các giá trị mặc định của RSEM.

2.2.3 Chú giải chức năng hệ phiên mã

2.2.3.1 Chú giải hệ phiên mã dựa trên cơ sở dữ liệu GO

Quá trình chú giải các unigene trên cơ sở dữ liệu GO được thực hiện bằng công cụ BLASTX trong phần mềm DIAMOND với ngưỡng E-value = $1.0E-5$. Quá trình chú giải sẽ đưa ra mã GO chứa đựng thông tin về trình tự và chức năng của các gen/ protein tham chiếu. Các mã chú giải GO được phân thành 3 nhóm chính dựa trên chức năng bao gồm các gen liên quan đến các quá trình sinh học (BP), thành phần tế bào (CC) và chức năng phân tử (MF). Trong 3 nhóm này, mỗi unigene sẽ được phân loại chi tiết hơn về nhóm chức năng.

2.2.3.2 Chú giải dựa trên cơ sở dữ liệu EggNOG

Chú giải các unigene trên cơ sở dữ liệu EggNOG bao gồm việc tìm kiếm, xác định các trình tự tương đồng trên Eukaryotic Orthologous Groups (KOGs), Cluster of Orthologous Groups (COGs) và Non-supervised Orthologous Groups (NOGs). Công cụ BLASTX trong phần mềm DIAMOND cũng được sử dụng để thực hiện chú giải các unigene dựa trên ngưỡng E-value là $1.0E-5$ [35]. Các unigene đã chú giải sẽ được sử dụng để tìm kiếm chú giải trên các nhóm tương đồng tương ứng trên cơ sở dữ liệu EggNOG.

2.2.3.3 Chú giải dựa trên cơ sở dữ liệu NT và NR của NCBI

Các unigene được chú giải sử dụng công cụ BLASTN của NCBI [36]. Giá trị ngưỡng E value được sử dụng trong quá trình chú giải là $1.0E-5$. Cơ sở dữ liệu NT cung cấp thông tin về trình tự các nucleic acid tham chiếu. Cơ sở dữ liệu NR chứa thông tin về trình tự các chuỗi amino acid tham chiếu.

2.2.3.4 Chú giải dựa trên cơ sở dữ liệu KEGG

Để chú giải chức năng của các unigene thu được trong nghiên cứu dựa trên cơ sở dữ liệu KEGG, công cụ BLASTX của phần mềm DIAMOND được sử dụng với ngưỡng giá trị E-value là $1.0E-5$. Việc tìm kiếm trình tự tương đồng được thực hiện bằng phương pháp Bi-directional Best Hit (BBH). Kết quả tìm kiếm và chú giải là các mã KO chứa các thông tin về protein và các con đường chuyển hóa, tổng hợp tham chiếu liên quan.

2.2.3.5 Chú giải dựa trên cơ sở dữ liệu UniProt

Quá trình chú giải các unigene trên cơ sở dữ liệu UniProt được thực hiện bằng công cụ BLASTX trong phần mềm DIAMOND với ngưỡng E-value = $1.0E-5$. Các trình tự unigene chú giải sẽ được tìm kiếm tương đồng trên hệ thống thông tin về chức năng của các protein tham chiếu UniProt KnowledgeBase (UniProtKB). Cơ sở dữ liệu này cũng cung cấp thông tin về tên, trình tự amino acid, dữ liệu phân loại hoặc tham khảo... thông qua mã UniProtKB.

2.2.3.6 Chú giải dựa trên cơ sở dữ liệu Pfam

Việc chú giải các unigene trên cơ sở dữ liệu Pfam được thực hiện bằng công cụ BLASTX trong phần mềm DIAMOND với ngưỡng E-value = $1.0E-5$. Chú giải trên Pfam được thực hiện thông qua so sánh sắp hàng trình tự và sử dụng các mô hình Markov ẩn (Hidden Markov models, HMMs). Quá trình chú giải sẽ cung cấp thông tin về các domain hoạt động, nhờ đó dự đoán chức năng của các protein.

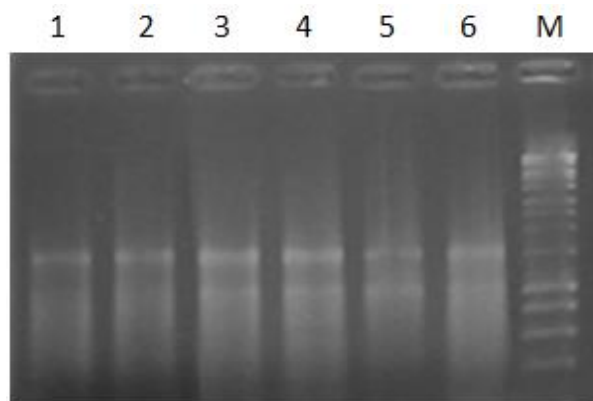
CHƯƠNG 3. KẾT QUẢ VÀ THẢO LUẬN

3.1 GIẢI TRÌNH TỰ HỆ PHIÊN MÃ ĐẶC HIỆU MÔ LÁ VÀ THÂN RỄ SÂM NGỌC LINH 4 NĂM TUỔI

3.1.1 Kết quả tách chiết, tinh sạch RNA tổng số và xây dựng các thư viện cDNA

RNA tổng số ở các mẫu lá và thân rễ của sâm Ngọc Linh được tách chiết, tinh sạch và đánh giá chất lượng bằng điện di trên gel agarose 8% (

Hình 3.1). Kết quả cho thấy, các mẫu RNA thu được có chất lượng tương đối tốt, có sự xuất hiện của các băng ribosome RNA với kích thước khoảng 1,5-2 kb. Để tinh sạch một số mẫu RNA, DNase được sử dụng nhằm loại bỏ DNA tổng số.



Hình 3.1 Kiểm tra kết quả tách chiết và tinh sạch RNA tổng số trên gel agarose

M: Thang chuẩn DNA 1 kb; trong đó, số thứ tự tương ứng với các mẫu: 1-3: C4.1-C4.3, 4-6: L4.1-L4.3

Kết quả kiểm tra nồng độ và độ sạch của RNA (Bảng 3.1) cho thấy, các sản phẩm có chỉ số A260/A280 thể hiện sự tinh sạch của mẫu dao động từ 1,90-2,15 chứng tỏ các mẫu RNA tách chiết được đủ điều kiện để tiến hành những thí nghiệm tiếp theo. Các RNA sau đó được phân thành những đoạn có kích thước

khoảng 450-550 bp sử dụng máy M220 Focused Ultrasonicator, đáp ứng các yêu cầu và được sử dụng để tổng hợp và xây dựng thư viện cDNA phục vụ giải trình tự gen thế hệ mới.

Bảng 3.1 Nồng độ và độ sạch (A260/A280) của các mẫu RNA sâm Ngọc Linh sau tách chiết và tinh sạch

STT	Tên mẫu	Nồng độ (ng/ μ l)	A260/A280
1	L4.1	1.010,0	1,96
2	L4.2	792,20	2,05
3	L4.3	1.468,3	2,08
4	C4.1	207,10	2,11
5	C4.2	482,40	2,11
6	C4.3	470,80	1,90

3.1.2 Kết quả giải trình tự hệ phiên mã của mô lá và mô thân rễ sâm Ngọc Linh 4 năm tuổi

Các thư viện được sử dụng cho giải trình tự thế hệ mới Illumina. Để đánh giá chất lượng của các đoạn đọc thô thu được sau khi giải trình tự, phần mềm FastQC được sử dụng để phân tích các chỉ số quan trọng của bộ dữ liệu như số lượng, kích thước của các đoạn đọc, % GC, điểm chất lượng trung bình (Q) của các trình tự theo từng base của các đoạn đọc, mức độ lặp của các đoạn đọc... Điểm chất lượng được biểu diễn dưới dạng biểu đồ hộp, trong đó hộp màu vàng thể hiện 50% số lượng phân bố các điểm chất lượng (gọi là khoảng interquartile), đường màu xanh chỉ giá trị trung bình. Điểm chất lượng trung bình 20 nghĩa là 99% độ chính xác và các đoạn đọc vượt qua điểm 20 được cho là chất lượng tốt. Điểm chất lượng trung bình từ 28 trở lên được đánh giá là chất lượng đọc rất tốt và đáng tin cậy. Điểm chất lượng trung bình bằng 30 tương đương với độ tin cậy 99,9% [37]. Thông tin thống kê cơ bản của kết quả giải trình tự các thư viện cDNA (Bảng 3.2) cho thấy đã thu được một số lượng lớn

các đoạn đọc từ các mẫu mô lá, thân rễ của sâm Ngọc Linh 4 năm tuổi với chất lượng khá cao, có độ tin cậy và đạt yêu cầu cho những phân tích tiếp theo.

Bảng 3.2 Thống kê các bộ dữ liệu thô khi giải trình tự các thư viện cDNA

Mẫu	Tổng số đoạn đọc	Tổng số base đọc	GC (%)	Q20 (%)	Q30 (%)
L4.1	89.888.850	9.078.773.850	43,33	98,89	96,25
L4.2	75.082.240	7.583.306.240	44,53	98,95	96,43
L4.3	66.477.494	6.714.226.894	43,15	98,20	94,39
C4.1	62.126.214	6.274.747.614	44,06	98,89	96,26
C4.2	78.869.048	7.925.373.848	42,70	97,14	92,04
C4.3	63.378.710	6.401.249.710	42,28	98,79	96,04
Tổng	435.822.556	43.977.678.156			

3.2 PHÂN TÍCH VÀ LẮP RÁP CÁC HỆ PHIÊN MÃ

3.2.1 Kết quả kiểm tra chất lượng các đoạn đọc sau khi trimming

Các bộ dữ liệu thô, sau khi được kiểm tra chất lượng bằng FastQC, tiếp tục được xử lý nhằm loại bỏ trình tự adapter và các base có điểm chất lượng thấp hơn ba từ hai đầu sử dụng Trimmomatic. Theo đó, các đoạn đọc có độ dài ngắn hơn 36 bp sẽ được loại bỏ để tạo ra dữ liệu sau khi trimming (Bảng 3.3). Kết quả, 429.930.834 các đoạn đọc tương ứng với 43.228.319.306 base đã thu được.

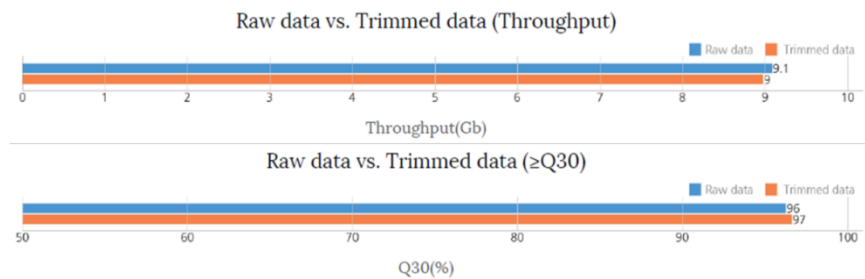
Bảng 3.3 Thống kê các bộ dữ liệu thu được sau khi trimming

Mẫu	Tổng số đoạn đọc	Tổng số base đọc	GC (%)	Q20 (%)	Q30 (%)
L4.1	89.141.174	8.975.651.153	43,34	99,12	96,61
L4.2	74.486.684	7.496.686.220	44,53	99,17	96,77
L4.3	65.523.252	6.583.474.260	43,17	98,64	95,01
C4.1	61.627.462	6.203.290.729	44,07	99,12	96,62

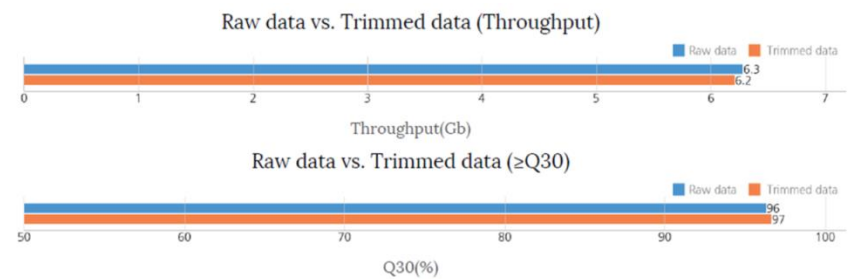
C4.2	76.350.630	7.648.769.112	42,74	97,88	93,14
C4.3	62.801.632	6.320.447.832	42,29	99,06	96,45
Tổng	429.930.834	43.228.319.306			

Sau khi trimming, điểm chất lượng trung bình của các trình tự theo từng base của các đoạn đọc cũng được đánh giá lại. So sánh chất lượng đoạn đọc của dữ liệu thô và dữ liệu sau khi trimming cho thấy dữ liệu giải trình tự thô ban đầu có chất lượng khá tốt và bước trimming để loại bỏ những phần không đạt yêu cầu giúp cho kết quả giải trình tự các đoạn đọc được tốt hơn (Hình 3.2). Kết quả lọc chất lượng giúp giảm số lượng đoạn đọc chất lượng thấp thông qua loại bỏ những vùng trình tự kém chất lượng trong các đoạn đọc. Điều này khiến cho số lượng đoạn đọc thu được sau bước lọc giảm nhưng chất lượng của các đoạn đọc lại tăng. Nhìn chung, phân tích điểm chất lượng trung bình của các đoạn đọc chỉ ra kết quả giải trình tự gen thể hệ mới với các mẫu cDNA của mô lá và mô thân rễ sâm Ngọc Linh thu được là tương đối tốt và đủ điều kiện cho lắp ráp và chú giải hệ phiên mã.

L4.1



C4.1



Hình 3.2 So sánh dữ liệu thô và dữ liệu sau trimming ở mẫu mô lá (L4.1) và thân rễ (C4.1) của sâm Ngọc Linh 4 năm tuổi

3.2.2 Kết quả lắp ráp *de novo* các hệ phiên mã

Quá trình phân tích dữ liệu trình tự hệ phiên mã của sâm Ngọc Linh bắt đầu bằng lắp ráp *de novo* các đoạn đọc trình tự đã chọn lọc chất lượng để tạo ra các contig là những đoạn trình tự có kích thước lớn hơn. Ở bước đầu tiên, phần mềm Trinity được sử dụng để chia nhỏ dữ liệu trình tự nhằm phân tích độc lập bằng biểu đồ de Bruijn tương ứng với các gen hay locus. Theo đó, các đoạn đọc được lắp ráp gói lên nhau để nối thành những phân đoạn dài hơn mà không chứa các gap hay N. Quá trình lắp ráp *de novo* được thực hiện thông qua 3 module phần mềm riêng biệt của Trinity là Inchworm, Chrysalis và Butterfly. Sau quá trình lắp ráp, phần mềm Trinity sẽ tạo ra một tệp “Trinity.fasta” chứa thông tin về trình tự phiên mã. Các đoạn trình tự thuộc cùng một locus sẽ được phân nhóm thành các cluster dựa trên độ tương đồng về trình tự. Những cluster phiên mã này được tạm coi là các gen.

Bảng 3.4 thể hiện kết quả thống kê của toàn bộ các contig được lắp ráp ban đầu từ các đoạn đọc thu được sau quá trình lọc chất lượng bao gồm số lượng gen, số lượng bản phiên mã (transcript), tỉ lệ GC, chỉ số N50, kích thước trung bình của contig và tổng số base được lắp ráp. Dữ liệu cho thấy kết quả lắp ráp hệ phiên mã của các mẫu mô lá và thân rễ sâm ở 4 năm tuổi có sự khác nhau về số lượng gen (47.870-70.526), số lượng bản phiên mã (69.638-106.276) và số lượng base lắp ráp (44.043.555-95.505.500). Cụ thể, các mẫu mô lá sâm Ngọc Linh 4 năm tuổi khác nhau về số lượng gen (68.454-81.480), số lượng bản phiên mã (99.609-121.554) và số lượng base lắp ráp (66.828.420-95.505.500); các mẫu thân rễ sâm Ngọc Linh 4 năm tuổi khác nhau về số lượng gen (47.870-63.268), số lượng bản phiên mã (69.638-100.651) và số lượng base lắp ráp (44.043.555-69.878.264). Kết quả lắp ráp chi tiết các đoạn đọc của các mẫu mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi được trình bày trên Bảng 3.5. Các chỉ số như số lượng gen, số lượng bản phiên mã, tỉ lệ GC, chỉ số N90-N10, kích thước các contig và số lượng base lắp ráp được thống kê cho toàn bộ các contig thu được và cho riêng contig có isoform với kích thước lớn nhất sau lắp ráp. Trong

đó, chỉ số N50 thể hiện cho ít nhất 50% các nucleotide được lắp ráp được tìm thấy trong các contig có chiều dài nhỏ nhất và cho isoform dài nhất tương ứng với các giá trị 993, 915, 1.074 của các mẫu mô lá sâm Ngọc Linh 4 năm tuổi; 861, 927, 816 của các mẫu thân rễ sâm Ngọc Linh 4 năm tuổi.

Bảng 3.4 Kết quả thống kê của contig được lắp ráp đầu tiên

Mẫu	Số lượng gen	Số lượng transcript	GC (%)	N50 (bp)	Kích thước trung bình của contig (bp)	Tổng số base được lắp ráp (bp)
L4.1	70.526	106.276	40,59	1.127	728,62	77.434.527
L4.2	68.454	99.609	40,89	1.033	670,91	66.828.420
L4.3	81.480	121.554	40,08	1.325	785,70	95.505.500
C4.1	47.870	69.638	40,20	870	632,46	44.043.555
C4.2	63.268	100.651	39,01	1.023	694,26	69.878.264
C4.3	59.717	88.228	38,86	828	613,73	54.148.380

Bảng 3.5 Kết quả thống kê quá trình lắp ráp các đoạn đọc ở mẫu mô lá (L4.1) và thân rễ (C4.1)

Mẫu	Thông số lắp ráp	Tất cả các contig của transcript	Chi isoform dài nhất/"gen"
L4.1	Tổng số trinity "gen"	70.526	70.526
	Tổng số trinity transcript	106.276	70.526
	GC (%)	40,59	40,80
	N90	300	258
	N80	458	347
	N70	665	481
	N60	888	699
	N50	1.127	993
	N40	1.374	1.312
	N30	1.645	1.642
	N20	1.989	2.025
	N10	2.547	2.618

	Độ dài contig lớn nhất	7.794	7.794
	Độ dài contig nhỏ nhất	201	201
	Median của độ dài contig	460,0	356,0
	Độ dài contig trung bình	728,62	623,76
	Tổng số base được lắp ráp	77.434.527	43.990.999
C4.1	Tổng số trinity “gen”	47.870	47.870
	Tổng số trinity transcript	69.638	47.870
	GC (%)	40,20	40,55
	N90	291	263
	N80	430	361
	N70	587	512
	N60	736	697
	N50	870	861
	N40	1.008	1.015
	N30	1.147	1.161
	N20	1.315	1.336
	N10	1.568	1.590
	Độ dài contig lớn nhất	7.854	7.854
	Độ dài contig nhỏ nhất	201	201
	Median của độ dài contig	494,0	402,0
	Độ dài contig trung bình	632,46	589,77
	Tổng số base được lắp ráp	44.043.555	28.232.335

3.2.3 Kết quả phân nhóm các đoạn trình tự thành các unigene

Các đoạn đọc sau khi lắp ráp *de novo* được tiến hành phân nhóm để tìm ra các unigene sử dụng CD-HIT-EST. Đây là một thuật toán bắt đầu với trình tự đầu vào có kích thước lớn nhất được chọn làm nhóm đại diện đầu tiên và phân chia những trình tự còn lại thành trình tự đại diện hay dư thừa dựa trên độ tương đồng với các đại diện đang xét. Độ tương đồng trình tự được tính toán dựa trên số lượng các word chung bằng cách sử dụng word indexing và bảng đếm để lọc ra những so sánh trình tự không cần thiết và tính độ tương đồng. Các bản phiên mã có kích thước lớn nhất thuộc cùng một locus sau khi được phân nhóm sẽ được coi là các contig unigene. Bảng 3.6 thống kê các contig unigene sau khi phân nhóm bao gồm số lượng gen, số lượng bản phiên mã, tỉ lệ GC, chỉ số N50,

kích thước trung bình của contig và tổng số base được lắp ráp. Các dữ liệu cho thấy, các mẫu mô khác nhau về số lượng gen (46.034-67.882), số lượng bản phiên mã (46.034-67.882) và số lượng base lắp ráp (27.641.662-50.440.507). Sự chênh lệch này chủ yếu do ảnh hưởng của số lượng đoạn đọc cũng như số lượng và kích thước các contig được lắp ráp.

Bảng 3.6 Kết quả thống kê của contig unigene

Mẫu	Số lượng gen	Số lượng transcript	GC (%)	N50 (bp)	Kích thước trung bình của contig (bp)	Tổng số base được lắp ráp (bp)
L4.1	67.882	67.882	40,73	1.022	635,07	43.109.582
L4.2	66.796	66.796	41,00	932	594,70	39.723.558
L4.3	77.381	77.381	39,97	1.120	651,85	50.440.507
C4.1	46.034	46.034	40,49	876	600,46	27.641.662
C4.2	59.818	59.818	38,90	961	615,62	36.825.282
C4.3	57.892	57.892	39,04	828	583,11	33.757.379

3.2.4 Kết quả dự đoán khung đọc mở

Bảng 3.7 thể hiện dữ liệu thống kê kết quả dự đoán ORF của các mẫu mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi. Kết quả thống kê cho thấy, đối với 3 mẫu L4.1, L4.2 và L4.3, có khoảng 33,04-37,01% các unigene được dự đoán ORF. Trong số đó, có 92,65-96,17% được dự đoán có 1 ORF và 3,83-7,35% có nhiều ORF. Trong số 22.091-25.713 ORF được dự đoán có 32,68%-48,4% ORF hoàn chỉnh, 25,45-34,38% ORF thiếu bộ ba mở đầu, 5,76-7,72% ORF thiếu bộ ba kết thúc và 18,71-27,18% ORF thiếu cả 2 loại này. Đối với 3 mẫu C4.1, C4.2 và C4.3, có khoảng 23,22-36,18% các unigene được dự đoán ORF. Trong số đó, có 95,61-97,95% được dự đoán có 1 ORF và 2,05-4,39% có nhiều ORF. Trong số 14.513-20.695 ORF được dự đoán có 22,22-36,53% ORF hoàn chỉnh, 46,65-

65,9% ORF thiếu bộ ba mở đầu, 1,73-4,09% ORF thiếu bộ ba kết thúc và 10,15-12,72% ORF thiếu cả 2 loại này. Kết quả, không có sự khác nhau nhiều giữa các mẫu dựa trên tỉ lệ tương đối giữa số lượng unigene và ORF.

Bảng 3.7 Kết quả thống kê dự đoán ORF của các mẫu mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi

Mẫu	Số lượng unigen	Unigen dự đoán là ORF	Unigen dự đoán là đơn ORF	Unigen dự đoán là đa ORF	Số lượng ORF	Gen hoàn chỉnh	Một phần/vùng mang mã	Một phần của đầu 5'	Một phần của đầu 3'
L4.1	67.882	22.426 (33,04%)	21.293 (94,95%)	1.133 (5,05%)	23.619	9.311 (39,42%)	5.659 (23,96%)	6.826 (28,9%)	1.823 (7,72%)
L4.2	66.796	24.722 (37,01%)	23.775 (96,17%)	947 (3,83%)	25.713	8.403 (32,68%)	6.989 (27,18%)	8.840 (34,38%)	1.481 (5,76%)
L4.3	77.381	20.458 (26,44%)	18.954 (92,65%)	1.504 (7,35%)	22.091	10.692 (48,4%)	4.134 (18,71%)	5.622 (25,45%)	1.643 (7,44%)
C4.1	46.034	16.653 (36,18%)	16.272 (97,71%)	381 (2,29%)	17.044	4.967 (29,14%)	1.939 (11,38%)	9.650 (56,62%)	488 (2,86%)
C4.2	59.818	13.887 (23,22%)	13.277 (95,61%)	610 (4,39%)	14.513	5.302 (36,53%)	1.846 (12,72%)	6.771 (46,65%)	594 (4,09%)
C4.3	57.892	20.270 (35,01%)	19.855 (97,95%)	415 (2,05%)	20.695	4.599 (22,22%)	2.101 (10,15%)	13.638 (65,9%)	357 (1,73%)

3.2.5 Kết quả ước lượng độ phong phú

Các đoạn đọc sau lọc chất lượng ở các mẫu mô lá và mô thân rễ sâm Ngọc Linh 4 năm tuổi được so sánh với trình tự tham chiếu của chính các mẫu này đã lắp ráp sử dụng phần mềm Bowtie. Để phân tích sự khác biệt biểu hiện gen, độ phong phú của các unigene giữa các mẫu được ước tính thông qua số lượng đoạn đọc sử dụng thuật toán RSEM. Bảng 3.8 thể hiện tỉ lệ lắp ráp các đoạn đọc của các mẫu mô lá và thân rễ sâm Ngọc Linh. Kết quả phân tích cho thấy, có 66,92-72,76% các đoạn đọc được lắp ráp vào hệ phiên mã tham chiếu

của chính các mẫu này. Điều này cho thấy có 27,24-33,08% các đoạn đọc không được lắp ráp vào hệ gen tham chiếu.

Bảng 3.8 Tỷ lệ lắp ráp các đoạn đọc của mẫu mô lá và thân rễ sâm

Mẫu	Số lượng đoạn đọc được xử lý	Số lượng đoạn đọc được lắp ráp	Số lượng đoạn đọc không được lắp ráp
L4.1	89.141.174	59.649.098 (66,92%)	29.942.076 (33,08%)
L4.2	74.486.684	50.129.670 (67,30%)	24.357.014 (32,70%)
L4.3	65.523.252	45.102.590 (68,83%)	20.420.662 (31,17%)
C4.1	61.627.462	44.405.868 (72,06%)	17.221.594 (27,94%)
C4.2	76.350.630	52.520.512 (68,79%)	23.830.118 (31,21%)
C4.3	62.801.632	45.693.642 (72,76%)	17.107.990 (27,24%)

3.3 CHÚ GIẢI HỆ PHIÊN MÃ CỦA CÁC MÔ SÂM NGỌC LINH 4 NĂM TUỔI

3.3.1 Chú giải hệ phiên mã mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi

Các unigene thu được ở các mẫu mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi được chú giải sử dụng các cơ sở dữ liệu Gene Ontology (GO), Evolutionary Genealogy of Genes: Non-supervised Orthologous Groups (EggNOG), NCBI Nucleotide (NT), NCBI non-redundant Protein (NR), Kyoto Encyclopedia of Genes and Genomes (KEGG), Universal Protein Resource (UniProt) và Pfam.

3.3.1.1 Kết quả chú giải dựa trên cơ sở dữ liệu GO

Kết quả chú giải các transcript của hệ phiên mã mẫu mô lá và thân rễ sâm Ngọc Linh dựa trên sự tương đồng về mặt trình tự trên cơ sở dữ liệu GO sẽ tương ứng với các mã GO chứa thông tin về trình tự và chức năng của các gen hoặc các sản phẩm của gen tham chiếu. Các transcript đã chú giải được phân chia theo chức năng thành 3 nhóm chính: Chu trình sinh học (biological process

- BP), thành phần tế bào (cellular component - CC) và chức năng phân tử (molecular function - MF). Trong nghiên cứu này, kết quả chú giải chức năng cho tỉ lệ các transcript thuộc nhóm BP chiếm 13,80-18,09% ở mô lá cây 4 năm tuổi và 14,05-17,53% ở mô thân rễ cây 4 năm tuổi. Các unigene thuộc nhóm CC chiếm từ 15,15-19,37% ở mô lá cây 4 năm tuổi và 15,75-19,50% ở mô thân rễ cây 4 năm tuổi. Tỉ lệ này của các transcript thuộc nhóm MF là 13,41-17,35% ở mô lá cây 4 năm tuổi và 14,11-17,50% ở mô thân rễ cây 4 năm tuổi. Số lượng unigene không được chú giải chiếm từ 45,19-57,64% ở mô lá cây 4 năm tuổi và 45,46-56,08% ở mô thân rễ cây 4 năm tuổi. Có thể thấy, tỉ lệ các unigene không được chú giải tương đối cao. Điều này có thể giải thích do sự hạn chế về số lượng các gen/ sản phẩm của gen tham chiếu cho các loài thuộc chi *Panax* trên cơ sở dữ liệu. Hình 3.3 minh họa thông tin phân nhóm của các mã GO thu được sau quá trình chú giải mẫu mô L4.1 sâm 4 năm tuổi.

3.3.1.2 Kết quả chú giải dựa trên cơ sở dữ liệu EggNOG

Thông tin về tỉ lệ phân nhóm dựa theo chức năng của các unigene đã chú giải ở mẫu mô lá L4.1 dựa trên cơ sở dữ liệu EggNOG được minh họa trên Hình 3.4. Các unigene được chú giải trên cơ sở dữ liệu EggNOG được chia làm 3 nhóm chính: (1) Nhóm các gen có chức năng trong các quá trình sinh học và lưu trữ thông tin (information storage and processing); (2) Nhóm các gen liên quan đến chu trình và tín hiệu tế bào (cellular processes and signaling); (3) Nhóm các gen liên quan đến trao đổi chất (metabolism). Tuy nhiên, tỉ lệ các unigene không được chú giải vẫn tương đối cao do hạn chế về thông tin của các loài thuộc chi *Panax* trên cơ sở dữ liệu EggNOG.

3.3.1.3 Kết quả chú giải dựa trên cơ sở dữ liệu NT và NR của NCBI

Dữ liệu về trình tự nucleotide và trình tự amino acid trên cơ sở dữ liệu NT và NR của NCBI được sử dụng làm tham chiếu để chú giải các hệ phiên mã. Quá trình tìm kiếm tương đồng các unigene của hệ phiên mã sâm Ngọc Linh đưa ra kết quả là các mã NCBI đại diện tương ứng với các unigene so sánh, chứa

thông tin về trình tự, tên chú giải, kích thước, độ bao phủ, tỉ lệ tương đồng, số lượng gap, độ tin cậy, vị trí khung đọc mở... (Bảng 3.).

3.3.1.4 Kết quả chú giải dựa trên cơ sở dữ liệu KEGG

Cơ sở dữ liệu KEGG chứa chủ yếu thông tin về trình tự, chức năng và các pathway liên quan của các protein tham chiếu. Bảng 3. trình bày kết quả chú giải hệ phiên mã của mẫu mô lá và thân rễ sâm Ngọc Linh đại diện. Hình 3.6 mô tả thông tin phân nhóm của các mã GO thu được sau quá trình chú giải các mẫu. Tương ứng với các unigene là các mã KEGG chứa thông tin về trình tự, tên chú giải, pathway và các thông tin so sánh khác như kích thước, độ bao phủ, tỉ lệ tương đồng, số lượng gap, độ tin cậy, vị trí khung đọc mở...

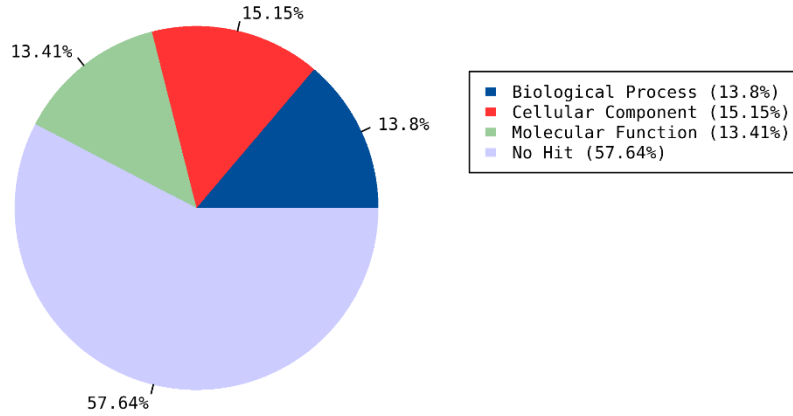
3.3.1.5 Kết quả chú giải dựa trên cơ sở dữ liệu UniProt

Dữ liệu về trình tự amino acid / protein trên cơ sở dữ liệu UniProt được sử dụng làm tham chiếu cho quá trình chú giải hệ phiên mã mẫu mô lá và thân rễ sâm Ngọc Linh. Bảng 3.1 trình bày kết quả chú giải hệ phiên mã đại diện. Tương ứng với các unigene là các mã UniProtKB chứa thông tin về trình tự và các thông tin so sánh khác như kích thước, độ bao phủ, tỉ lệ tương đồng, số lượng gap, độ tin cậy, vị trí khung đọc mở...

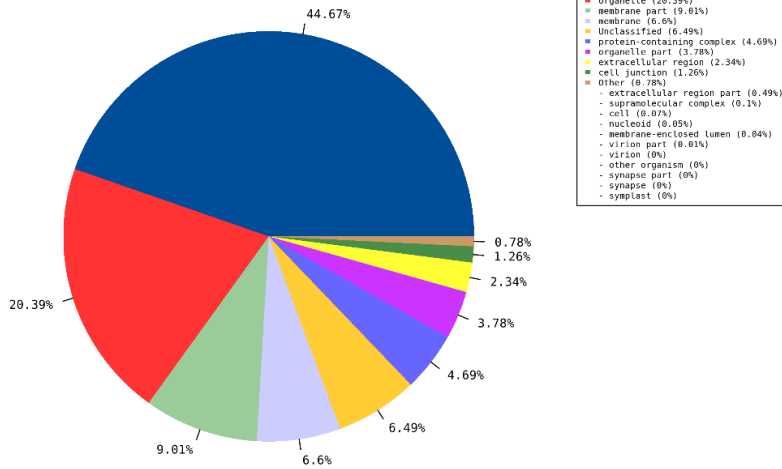
3.3.1.6 Kết quả chú giải dựa trên cơ sở dữ liệu Pfam

Cơ sở dữ liệu Pfam chứa thông tin về trình tự, chức năng và thông tin về các domain của các protein tham chiếu. Quá trình chú giải các unigene của hệ phiên mã mẫu mô lá và thân rễ sâm Ngọc Linh đưa ra kết quả là các mã chú giải Pfam tương ứng với các unigene so sánh, chứa thông tin về trình tự, tên chú giải và các thông tin so sánh khác như kích thước, độ bao phủ, tỉ lệ tương đồng, số lượng gap, độ tin cậy, vị trí khung đọc mở... (Bảng 3.11).

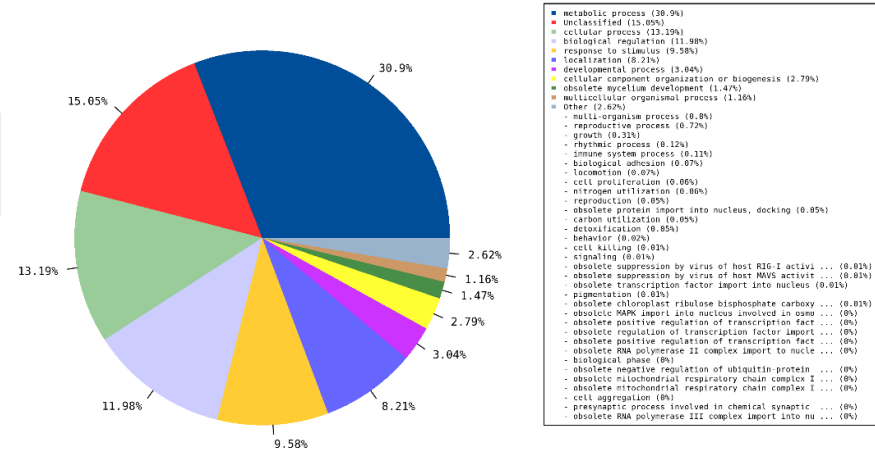
L4.1 GO Analysis Result



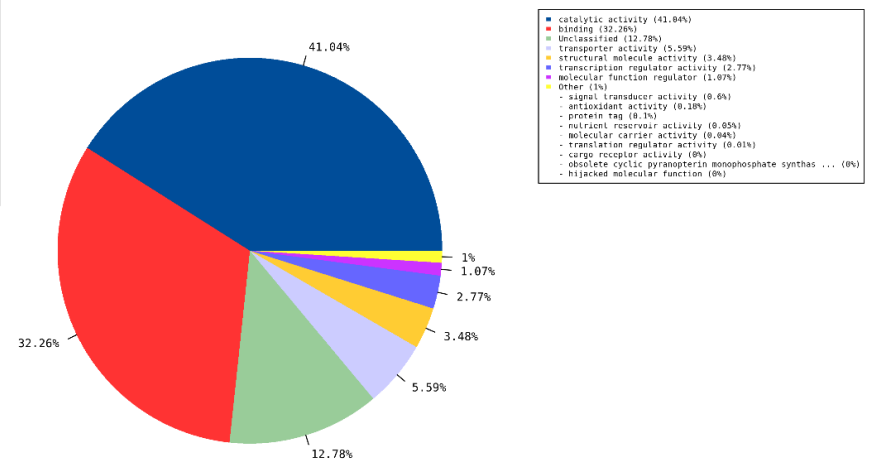
Cellular Component



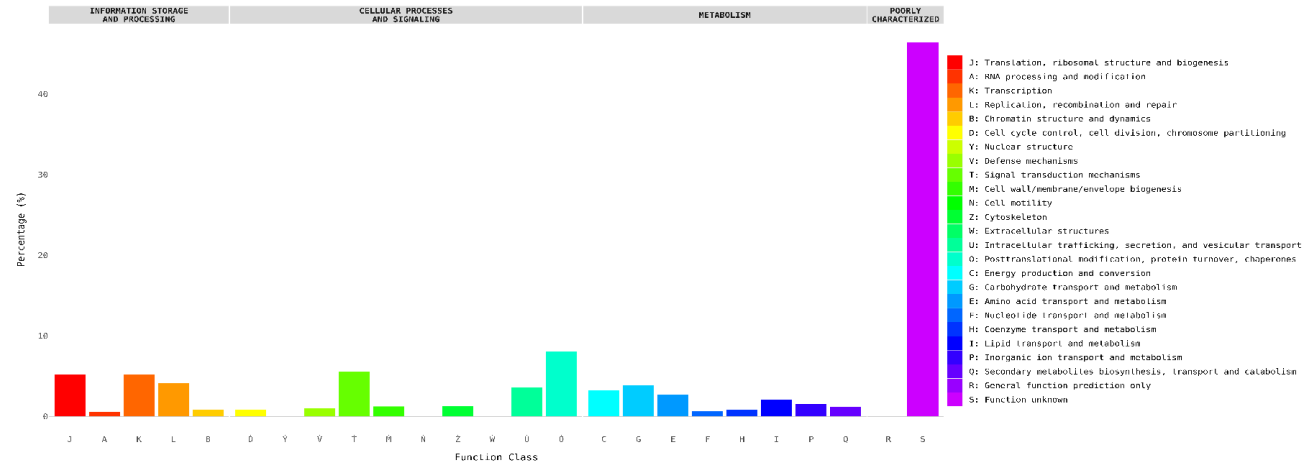
Biological Process



Molecular Function



L4.1



Hình 3.4. Kết quả chú giải và phân nhóm các gen thuộc hệ phiên mã của mẫu mô lá (L4.1) dựa trên cơ sở dữ liệu EggNOG

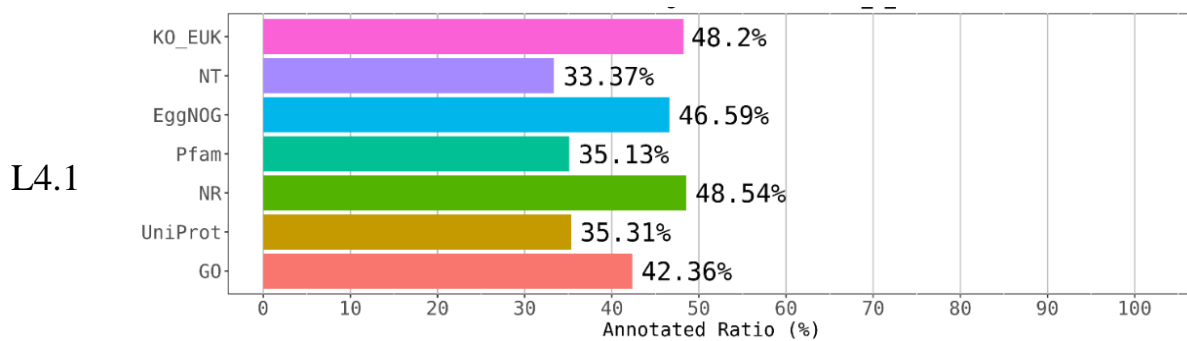
3.3.2 Kết quả chú giải các unigene của mẫu mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi

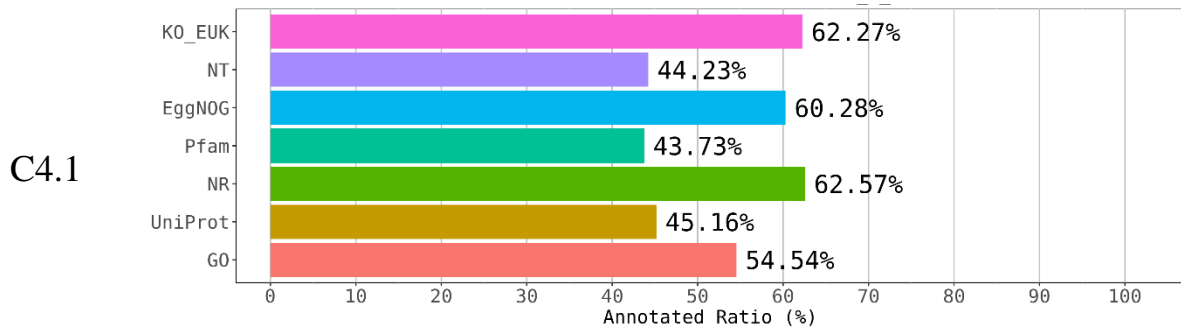
Quá trình chú giải hệ phiên mã mẫu mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi được thống kê trong Bảng 3.13. Kết quả tổng hợp cho thấy tỉ lệ các unigene được chú giải trên tổng số các unigene thu được chiếm 51,72-64,81% ở mẫu mô lá 4 tuổi và 52,87-66,32% ở mẫu thân rễ 4 tuổi.

Bảng 3.13 Tỉ lệ các unigene được chú giải trong hệ phiên mã các mẫu mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi

Mẫu	Số lượng unigene	Số lượng unigene được chú giải	Tỉ lệ chú giải (%)
L4.1	77.381	40.021	51,72
L4.2	66.796	43.289	64,81
L4.3	61.600	36.847	59,82
C4.1	46.034	30.528	66,32
C4.2	57.892	34.218	59,11
C4.3	59.818	31.624	52,87

Thông tin về tỉ lệ các unigene được chú giải ở từng mẫu trên từng cơ sở dữ liệu cũng được trình bày trong Hình 3.5. So sánh tỉ lệ chú giải thành công giữa các cơ sở dữ liệu nhận thấy KEGG, EggNOG và NR luôn đạt tỉ lệ chú giải thành công cao nhất, sau đó là GO, UniProt, Pfam và thấp nhất là NT.





Hình 3.5 Tỷ lệ các unigene của hệ phiên mã mẫu mô lá (L4.1) và thân rễ (C4.1) được chú giải trên các cơ sở dữ liệu

3.3.3 Tổng hợp, phân tích và so sánh dữ liệu hệ phiên mã ở mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi

Dữ liệu hệ phiên mã các mẫu đơn được chuẩn hóa để có thể so sánh đối chiếu. Kết quả so sánh tỷ lệ chú giải cụ thể của unigene sâm Ngọc Linh 4 năm tuổi dựa trên các cơ sở dữ liệu khác nhau được tổng hợp trên

Bảng 3.4. Số lượng unigene được chú giải chức năng bởi ít nhất một cơ sở dữ liệu trong số các cơ sở dữ liệu GO, EggNOG, NT, NR, KEGG, UniProt và Pfam ở mẫu sâm 4 năm tuổi là 60,85%. Tỷ lệ chú giải cao nhất cũng thuộc về cơ sở dữ liệu KEGG với tỷ lệ 57,95%. Trong khi đó, cơ sở dữ liệu có khả năng chú giải số unigene thấp nhất thuộc về Uniprot với 45,40%.

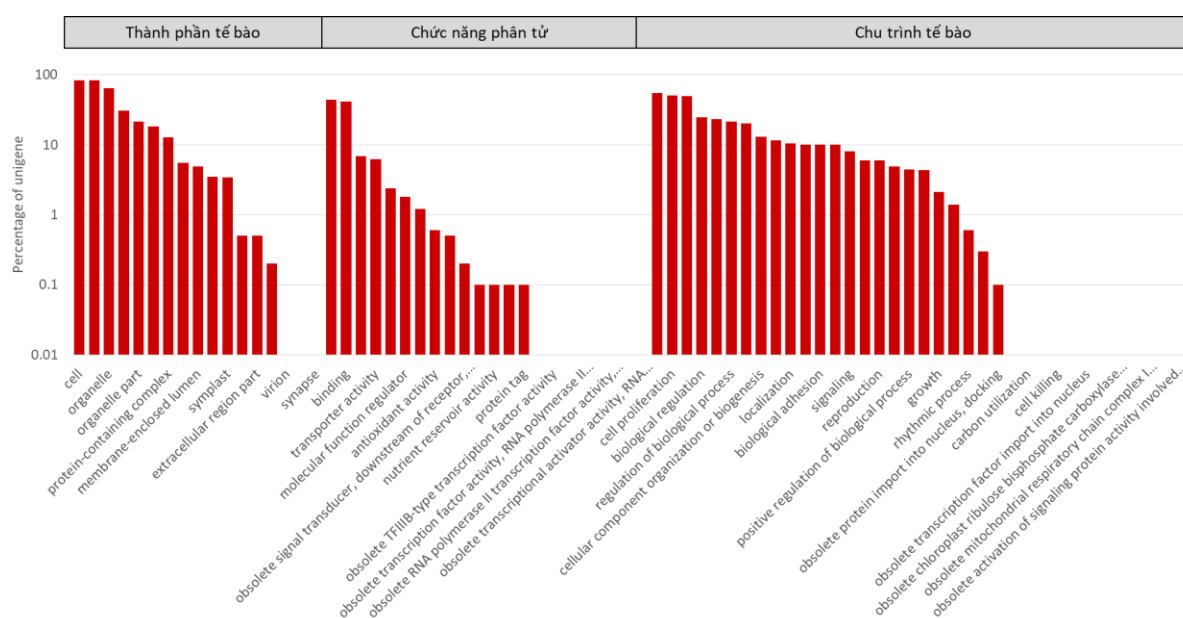
Bảng 3.14 Khả năng chú giải các unigene của hệ phiên mã sâm Ngọc Linh 4 năm tuổi trên các cơ sở dữ liệu khác nhau

Cơ sở dữ liệu	GO	Uniprot	NR	EggNOG	KEGG	NT	Pfam	Tổng
Số unigene	22.900	19.510	24.816	24.305	24.901	19.690	19.529	26.149
Tỷ lệ chú giải (%)	53,29	45,40	57,75	56,56	57,95	45,82	45,45	60,85

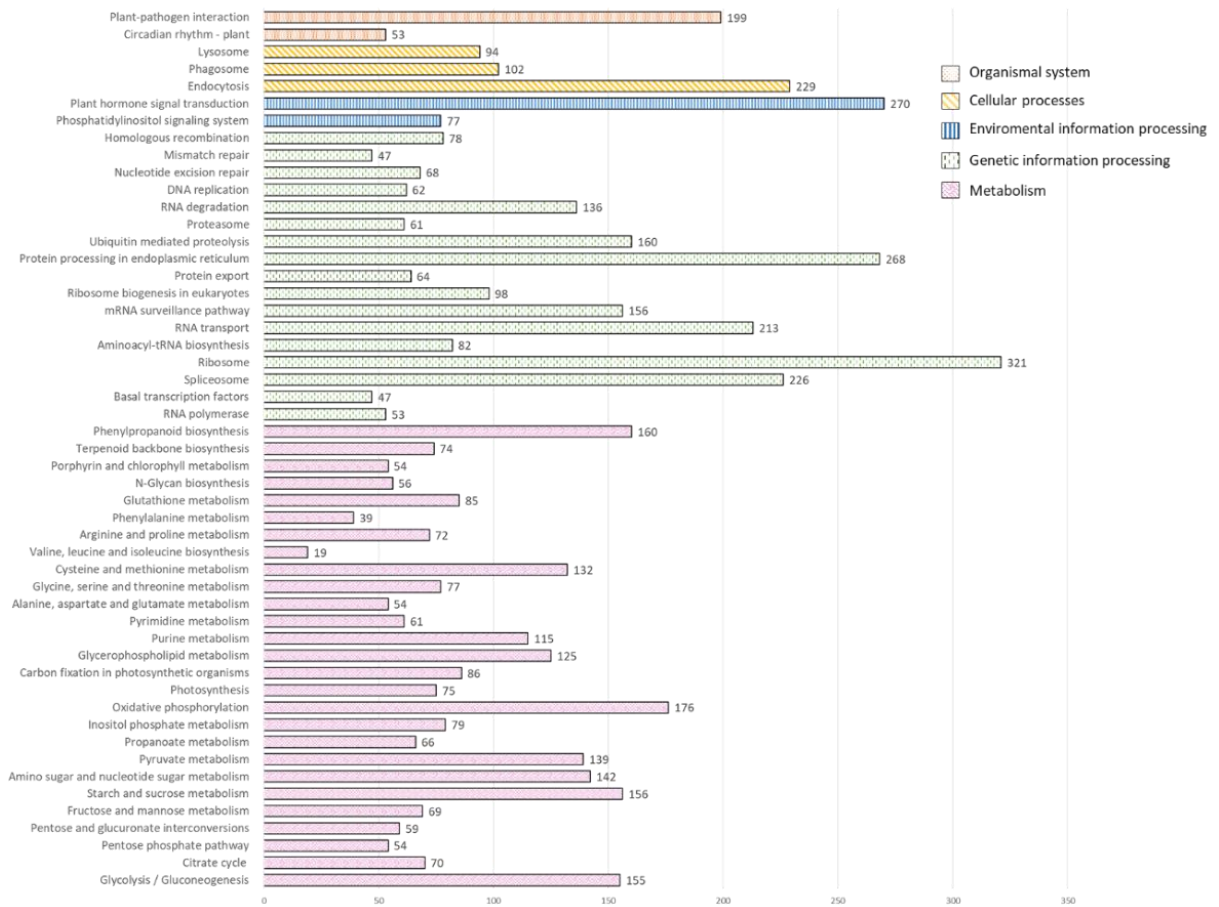
Với cơ sở dữ liệu GO, tất cả unigene được phân loại dựa trên 3 nhóm chức năng lớn và các nhóm nhỏ hơn. Đối với mẫu sâm Ngọc Linh 4 năm tuổi, tổng số 22.900 unigene được chú giải với 3 nhóm lớn CC, BP và MF có tỷ lệ lần

lượt là 47,81%, 36,96% và 15,23%. Trong nhóm CC, các unigene thuộc phân nhóm “tế bào/thành phần tế bào” chiếm tỉ lệ cao nhất với hơn 83%, tiếp đến là “bào quan” và “màng tế bào”. Trong nhóm BP, chủ yếu là các unigene mã cho protein tham gia “chu trình tế bào” và “quá trình trao đổi chất” với tỉ lệ lần lượt là 54,6% và 49,0%. Phần lớn unigene nằm trong nhóm MF được phân loại vào phân nhóm “hoạt động xúc tác” với 10.096 unigene và “liên kết” với 9.522 unigene (Hình 3.6).

Với cơ sở dữ liệu KEGG, chức năng unigene liên quan đến các con đường chuyển hóa trong sinh vật được đánh giá. Ở hệ phiên mã sâm 4 năm tuổi, 24.901 unigene đã được chú giải và chia vào 211 con đường chuyển hóa trong sinh vật (Hình 3.7). Các unigene được chú giải này sẽ là cơ sở cho các nghiên cứu tiếp theo về các hợp chất thứ cấp quan trọng ở cây sâm Ngọc Linh Việt Nam.

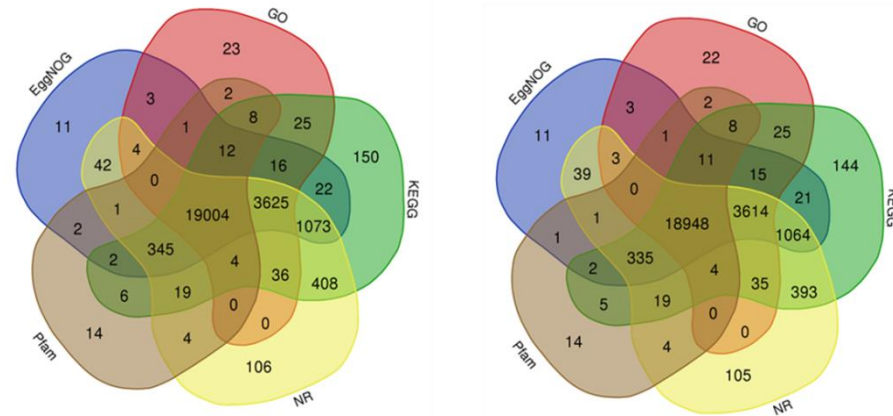


Hình 3.6 Chú giải GO cho unigene sâm Ngọc Linh 4 năm tuổi



Hình 3.7 Chú giải KEGG của hệ phiên mã sâm Ngọc Linh 4 năm tuổi

Hệ phiên mã của mô lá và mô thân rễ sâm Ngọc Linh 4 năm tuổi được so sánh dựa trên kết quả chú giải bởi các cơ sở dữ liệu GO, EggNOG, NR, KEGG và Pfam (Hình 3.8). Cơ sở dữ liệu có khả năng chú giải cao nhất là KEGG và cơ sở dữ liệu có khả năng chú giải thấp nhất trong 5 cơ sở dữ liệu trên là EggNOG.



Hình 3.8 Biểu đồ venn biểu diễn số unigene được chú giải bởi các cơ sở dữ liệu của hệ phiên mã mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi

Tiếp theo, sự có mặt của các transcript mang ưu thế ở hệ phiên mã mỗi loại mô sâm 4 năm tuổi được phân tích. Bảng 3.15 cho thấy một số transcript mã hóa cho các protein tác động tới khả năng chống chịu bệnh, chống stress hay sinh trưởng và già hóa đều được tìm thấy ở hai hệ phiên mã sâm 4 năm tuổi. Tuy nhiên, sự biểu hiện của chúng khi được quan sát thấy ở lá thấp hơn ở thân rễ (dựa vào số lượng read count). Gen biểu hiện mạnh nhất ở lá là “Senescence-associated protein” và cao thứ hai ở rễ. Điều này có thể do thời điểm thu mẫu vào đầu thu, thời gian bắt đầu chuyển sang mùa rụng lá. Một số gen mã cho protein liên quan đến yếu tố gây bệnh, catalase_like hay metallothionein_like được tăng cường biểu hiện vào thời gian này [38], cũng được thấy biểu hiện cao ở lá. Ngoài ra, gen mã cho protein liên quan lục lạp (chloroplast) có sự biểu hiện cao ở lá - loại mô đặc trưng chứa lục lạp.

Bảng 3.15 Danh sách transcript chiếm ưu thế trong hệ phiên mã của mô lá và thân rễ sâm Ngọc Linh 4 năm tuổi

TT	Unigene ID	Accession no.	Chú giải	Read count
Mô lá				
1	c698379_g1_i1	mtr:MTR_0055s0050	Senescence-associated protein	1.600.684
2	c801000_g7_i2	MF377623	<i>P. vietnamensis</i> chloroplast	420.495
3	c494428_g1_i1	DQ384524	<i>P. ginseng</i> GBR5-like protein mRNA	393.661
4	c242479_g3_i1	dcr:108218502	Catalase-like	201.600
5	c801000_g2_i1	MF377623	<i>P. vietnamensis</i> chloroplast	190.935
6	c854330_g1_i1	dcr:108220108	Metallothionein-like protein type 3	150.255
7	c795392_g4_i3	dcr:108214401	Heat shock cognate 70 kDa protein 2-like	138.584
8	c492933_g1_i1	KJ741402	<i>P. notoginseng</i> pathogenesis-related protein (PR10-1) mRNA	110.463
9	c801571_g4_i1	DQ384525	<i>P. ginseng</i> ADP-ribosylation factor-like protein mRNA	86.711
10	c789616_g1_i1	dcr:108212952	Uncharacterized protein	73.169
Mô thân rễ				
1	c797413_g2_i1	KC751542	<i>P. notoginseng</i> RNase-like major storage protein	2.195.014
2	c698379_g1_i1	mtr:MTR_0055s0050	Senescence-associated protein	870.817
3	c801571_g4_i1	DQ384525	<i>P. ginseng</i> ADP-ribosylation factor-like protein	270.295
4	c795392_g4_i3	dcr:108214401	Heat shock cognate 70 kDa protein 2-like	113.164
5	c800014_g1_i2	dcr:108218481	Delta(12)-fatty-acid desaturase FAD2-like	97.079
6	c840150_g1_i1	qsu:112036655	Pathogenesis-related protein 1-like	94.124

7	c791685_g1_i2	pop:7494466	18.1 kDa class I heat shock protein	83.537
8	c696765_g1_i1	dcr:108224592	Protein translation factor SUI1 homolog	77.054
9	c739909_g1_i1	dcr:108218482	Delta(12)-fatty-acid desaturase FAD2-like	75.405
10	c494485_g1_i1	dcr:108192288	22.0 kDa heat shock protein-like	71.853

Kết quả từ bảng trên cũng cho thấy, một số unigene mã cho protein tham gia quá trình già hóa, protein heat shock và protein liên quan yếu tố gây bệnh đều xuất hiện nhiều ở hệ phiên mã mô lá và mô thân rễ. Về hệ phiên mã mô thân rễ sâm 4 năm tuổi, đứng đầu trong danh sách các transcript chiếm ưu thế là unigene mã cho “RNase-like major storage protein”. Ngoài một số unigene mã cho protein heat shock hay liên quan yếu tố gây bệnh tương tự ở lá, có thể quan sát được một số unigene khác mã cho protein “ADP-ribosylation factor-like protein”, “Delta(12)-fatty-acid desaturase FAD2-like” và “Protein translation factor SUI1 homolog” biểu hiện cao ở mô thân rễ. “Delta(12)-fatty-acid desaturase FAD2-like” có hai đại diện unigene (vị trí 5 và 9). Họ gen *FAD2-like* mã cho các enzyme bao gồm hydrolase, epoxygenase, conjugase và acetylenase góp phần đa dạng hóa các loại chất béo [39].

Tóm lại, sử dụng các cơ sở dữ liệu GO, EggNOG, NT, NR, KEGG, UniProt và Pfam, các unigene thu được ở các mẫu mô thân rễ và lá của sâm Ngọc Linh 4 năm tuổi đã được chú giải với tỉ lệ đạt từ 60,88-66,29%. KEGG, EggNOG và NR có tỉ lệ chú giải thành công cao nhất, tiếp đó là GO, UniProt, Pfam và thấp nhất là NT. Với GO, các unigene được phân vào 3 nhóm chức năng lớn và 52 nhóm nhỏ. Với NR, các trình tự gen sâm tương đồng với trình tự của 534 loài sinh vật. Với KEGG, 24.901 unigene đã được chú giải và chia vào 211 con đường chuyển hóa trong sinh vật; trong đó, 74 unigene thuộc con đường sinh tổng hợp khung terpenoid đã được chú giải, là cơ sở cho các nghiên cứu tiếp theo về các hoạt chất thứ cấp mang giá trị cao ở cây sâm Ngọc Linh.

Ngoài ra, các unigene thuộc các hệ phiên mã mô thân rễ và lá sâm có sự tương đồng về tỉ lệ giữa các nhóm chức năng. Ở hệ phiên mã, các transcript có ưu thế thường liên quan tới các protein chống chịu bệnh, sinh trưởng, già hóa và chống stress.

Việc có được cái nhìn bao quát về mức độ biểu hiện của sâm Ngọc Linh ở các thời điểm, lứa tuổi khác nhau (cụ thể trong nghiên cứu này là 4 năm tuổi), hoặc các mô khác nhau (ở đây là mô lá và mô thân rễ), cho phép cân nhắc để tối ưu các yếu tố sinh trưởng, canh tác, thời điểm thu hoạch..., qua đó có thể thu được hàm lượng các hợp chất thứ cấp quan trọng cao nhất. Do đó, kết quả của nghiên cứu này là mảnh ghép quan trọng để có thể hiểu rõ hơn các con đường tổng hợp, mức độ biểu hiện khác biệt ở hai nhóm mô khác nhau ở 4 năm tuổi, cũng như là tiền đề cho các nghiên cứu tiếp theo ở lứa tuổi/nhóm mô khác; bên cạnh đó, có định hướng nuôi trồng, chăm sóc và thu hoạch sâm vào thời điểm phù hợp và có hiệu quả cao nhất. Các kết quả thu được cũng góp phần xác định quyền sở hữu quốc gia về nguồn gen đặc hữu, từ đó có thể khai thác để đánh giá và nâng cao chất lượng nguồn gen sâm Ngọc Linh Việt Nam thông qua công nghệ sinh học.

CHƯƠNG 4. KẾT LUẬN VÀ KIẾN NGHỊ

4.1 KẾT LUẬN

- Các hệ phiên mã của các mẫu mô thân rễ và lá sâm Ngọc Linh 4 năm tuổi được lắp ráp từ 429.930.834 các đoạn đọc, tương ứng với 43.228.319.306 base sau khi lọc chất lượng. Hệ phiên mã của các mẫu mô khác nhau về số lượng gen, số lượng bản phiên mã và số lượng base lắp ráp.

- Các bộ dữ liệu hệ phiên mã được phân nhóm thành các unigene, tìm kiếm, dự đoán ORF và ước lượng độ phong phú. Đối với hệ phiên mã sâm Ngọc Linh 4 năm tuổi, khoảng 23,22-37,01% các unigene được dự đoán là ORF. Tổng 66,92-72,76% các đoạn đọc của hệ phiên mã mẫu mô lá và thân rễ sâm Ngọc Linh được lắp ráp vào hệ phiên mã tham chiếu.

- Các unigene được chú giải với tỉ lệ đạt từ 51,72-64,81%, sử dụng các cơ sở dữ liệu GO, EggNOG, NT, NR, KEGG, UniProt và Pfam. KEGG, EggNOG và NR có tỉ lệ chú giải thành công cao nhất, tiếp đó là GO, UniProt, Pfam và thấp nhất là NT. Với GO, các unigene được phân chia vào 3 nhóm chức năng lớn và 52 nhóm nhỏ. Với KEGG, các unigene được chia vào 211 con đường chuyển hóa. Với NR, các trình tự gen sâm Ngọc Linh tương đồng với trình tự của 534 loài sinh vật.

- Các unigene thuộc các hệ phiên mã mô lá và thân rễ sâm có sự tương đồng về tỉ lệ giữa các nhóm chức năng. Các đoạn transcript mang ưu thế ở các hệ phiên mã có liên quan đến các protein chống stress, chống chịu bệnh, già hóa hay sinh trưởng.

4.2 KIẾN NGHỊ

Dữ liệu thu được từ nghiên cứu này là cơ sở thực tiễn để tiếp tục mở rộng các nghiên cứu xây dựng dữ liệu toàn diện về hệ phiên mã của sâm Ngọc

Linh ở các độ tuổi, góp phần cung cấp thông tin đầy đủ và hữu ích về hệ phiên mã; xa hơn là các gen chức năng liên quan đến các tính trạng quan trọng, phục vụ công tác đánh giá, so sánh nguồn gen quý của sâm Ngọc Linh ở các lứa tuổi khác nhau; hỗ trợ bảo tồn, chọn tạo, quản lý và nhân giống sâm Ngọc Linh bằng công nghệ sinh học.

TÀI LIỆU THAM KHẢO

1. Baeg IH, So SH (2013). The world ginseng market and the ginseng (Korea). *J Ginseng Res* 37(1): 1.
2. Yang MS, Wu MY (2016). Chinese ginseng. *Nutraceuticals* 50: 693-705.
3. Jung J, Lee NK, Paik HD (2017). Bioconversion, health benefits, and application of ginseng and red ginseng in dairy products. *Food Sci Biotechnol* 26: 1155-1168.
4. Patel S, Rauf A (2017). Adaptogenic herb ginseng (*Panax*) as medical food: status quo and future prospects. *Biomed Pharmacother* 85: 120-127.
5. Shin BK, Kwon SW, Park JH (2015). Chemical diversity of ginseng saponins from *Panax ginseng*. *J Ginseng Res* 39(4): 287-298.
6. Sách Đỏ Việt Nam. Phần II. Thực vật (2007). *Nhà xuất bản Khoa học tự nhiên và Công nghệ, Hà Nội*.
7. Ha TD, Grushvitzky IV (1985). A new species of the genus *Panax* (Araliaceae) from Vietnam. *Bot Zhurn* 70: 519-522.
8. Nguyễn Tập (2005). Các loài thuộc chi *Panax* ở Việt Nam. *Tạp chí Dược liệu* 10(3): 71-76.
9. Trần Ngọc Lâm, Nguyễn Tiến Dũng, Nguyễn Thị Thu, Lê Thị Thu Hiền, Ngô Hoàng Linh, Nguyễn Đức Nam, Trần Quốc Thành, Hoàng Nghĩa Nhạc, Phùng Văn Hào (2016). Kết quả nghiên cứu về loài sâm Puxailaileng ở vùng núi cao tỉnh Nghệ An. *Tạp chí Khoa học-Công nghệ Nghệ An* 12: 7-11.
10. Lương Đức Toàn (2018). Tính chất, chất lượng đặc thù của sâm củ và điều kiện tự nhiên vùng trồng sâm Ngọc Linh tỉnh Kon Tum. Kỷ yếu Hội nghị Đầu tư và phát triển sâm Ngọc Linh Kon Tum và các dược liệu

khác. Bộ Y tế, UBND Tỉnh Kon Tum, Bộ Nông nghiệp và Phát triển Nông thôn. Kon Tum.

11. Shendure J, Ji H (2008). Next-generation DNA sequencing. *Nat Biotechnol* 26: 1135-1145.
12. Metzker ML (2010). Sequencing technologies - the next generation. *Nat Rev Genet* 11: 31-46.
13. Liu L, Li Y, Li S, Hu N, He Y, Pong R, Lin D, Lu L, Law M (2012). Comparison of next-generation sequencing systems. *J Biomed Biotechnol*: 1-11.
14. Quail MA, Smith M, Coupland P, Otto TD, Harris SR, Connor TR (2012). A tale of three next generation sequencing platforms: Comparison of Ion Torrent, Pacific Biosciences and Illumina MiSeq sequencers. *BMC Genomics* 13: 341.
15. Ferrarini M, Cestaro A, Sargent DJ, Moretto M, Ward JA, Šurbanovski N, Stevanović V, Giongo L, Viola R, Cavalieri D, Velasco R, Cestaro A, Sargent DJ (2013). An evaluation of the PacBio RS platform for sequencing and *de novo* assembly of a chloroplast genome. *BMC Genomics* 14: 670-670.
16. Schuster SC (2008). Next-generation sequencing transforms today's biology. *Nat Methods* 5(1): 16-18.
17. Rusk N (2011). Torrents of sequence. *Nat Methods* 8: 44.
18. Poehlmann A, Kuester D, Meyer F, Lippert H, Roessner A, Schneider-Stock R (2007). Kras mutation detection in colorectal cancer using the Pyrosequencing technique. *Pathol Res Pract* 203: 489-497.
19. Pettersson E, Lundeberg J, Ahmadian A (2009). Generations of sequencing technologies. *Genomics* 93 (2): 105-111.
20. Schadt EE, Turner S, Kasarskis A (2010). Window into third-generation sequencing. *Hum Mol Genet* 19 (R2): R227-40.

21. Sun C, Li Y, Wu Q, Luo H, Sun Y, Song J, Lui EM, Chen S (2010). *De novo* sequencing and analysis of the *American ginseng* root transcriptome using a GS FLX Titanium platform to discover putative genes involved in ginsenoside biosynthesis. *BMC Genomics* 11: 262.
22. Luo H, Sun C, Sun Y, Wu Q, Li Y, Song J, Niu Y, Cheng X, Xu H, Li C, Liu J, Steinmetz A, Chen S (2011). Analysis of the transcriptome of *Panax notoginseng* root uncovers putative triterpene saponin-biosynthetic genes and genetic markers. *BMC Genomics* 5: S5.
23. Li B, Dewey CN (2011). RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics* 12: 323.
24. Cao H, Nuruzzaman M, Xiu H, Huang J, Wu K, Chen X, Li J, Wang L, Jeong JH, Park SJ, Yang F, Luo J, Luo Z (2015). Transcriptome analysis of methyl jasmonate-elicited *Panax ginseng* adventitious roots to discover putative ginsenoside biosynthesis and transport genes. *Int J Mol Sci* 16(2): 3035-3057.
25. Jayakodi M, Lee SC, Lee YS, Park HS, Kim NH, Jang W, Yang TJ (2015). Comprehensive analysis of *Panax ginseng* root transcriptomes. *BMC Plant Biol* 15(1): 1-12.
26. Rai A, Yamazaki M, Takahashi H, Nakamura M, Kojoma M, Suzuki H, Saito K (2016). RNA-seq transcriptome analysis of *Panax japonicus*, and its comparison with other *Panax* species to identify potential genes involved in the saponins biosynthesis. *Front Plant Sci* 7: 481.
27. Vu DD, Shah SNM, Pham MP, Bui VT, Nguyen MT, Nguyen TPT (2020). *De novo* assembly and transcriptome characterization of an endemic species of Vietnam, *Panax vietnamensis* Ha et Grushv., including the development of EST-SSR markers for population genetics. *BMC Plant Biol* 20: 358.

28. Andrews S, Krueger F, Segonds-Pichon A, Biggins L, Krueger C, Wignett S (2012). FastQC. A quality control tool for high throughput sequence data. Babraham Institute, Cambridge, United Kingdom.
29. Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: A flexible trimmer for Illumina Sequence Data. *Bioinformatics*, btu170.
30. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury R, Zeng Q, Chen Z, Mauceli E, Hacohen N, Gnirke A, Rhind N, di Palma F, Birren BW, Nusbaum C, Lindblad-Toh K, Friedman N, Regev A (2011). Full-length transcriptome assembly from RNA-seq data without a reference genome. *Nat Biotechnol* 29(7): 644-652.
31. Li W, Godzik A (2006). Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* 22(13): 1658-1659.
32. Fu L, Niu B, Zhu Z, Wu S, Li W (2012). CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics* 28(23): 3150-3152.
33. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D, Li B, Lieber M, Macmanes MD, Ott M, Orvis J, Pochet N, Strozzi F, Weeks N, Westerman R, William T, Dewey CN, Henschel R, Leduc RD, Friedman N, Regev A (2013). *De novo* transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. *Nat Protoc* 8(8): 1494-512.
34. Buchfink B, Xie C, Huson DH. (2015). Fast and sensitive protein alignment using DIAMOND. *Nat Methods* 12(1): 59-60.
35. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990). Basic local alignment search tool. *J Mol Biol* 215(3): 403-410.
36. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Pachter L (2010). Transcript assembly and abundance estimation

from RNA-Seq reveals thousands of new transcripts and switching among isoforms. *Nat Biotechnol* 28(5): 511.

37. Cock PJ, Fields CJ, Goto N, Heuer ML, Rice PM (2010). The Sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants. *Nucleic Acids Res* 38: 1767-1771.
38. Gepstein S, Sabehi G, Carp MJ, Hajouj T, Neshher MFO, Yariv I, Bassani M (2003). Large-scale identification of leaf senescence-associated genes. *Plant J* 36(5): 629-642.
39. Dyer JM, Chapital DC, Kuan JCW, Mullen RT, Turner C, McKeon TA, Pepperman AB (2002). Molecular analysis of a bifunctional fatty acid conjugase/desaturase from tung. Implications for the evolution of plant fatty acid diversity. *Plant Physiol* 130(4): 2027-2038.